

Robuste Verfahren für strukturierte
hochdimensionale
Repeated-Measures-Designs unter
Nicht-Normalverteilung

Diplomarbeit

vorgelegt von
Hans-Joachim Helms
aus Celle

angefertigt am
Institut für Mathematische Stochastik
der Georg-August-Universität Göttingen

2010

Danksagung

Zuerst möchte ich Herrn Prof. Dr. Brunner für die hervorragende Betreuung meiner Diplomarbeit danken. Außerdem danke ich Herrn Prof. Dr. Schlather für die Übernahme des Koreferats. Ebenfalls gilt mein Dank Herrn Dr. Konietschke für die Betreuung meiner Diplomarbeit.

Weiterhin möchte ich mich bei der gesamten Abteilung der medizinischen Statistik für die schöne Arbeitsatmosphäre, die aufheiternden Gespräche und den leckeren Kuchen bedanken.

Mein spezieller Dank gilt meinen Freunden, den Mitgliedern der Schweden-Wanderung, dem Sport und Mandy, die mich immer unterstützt hat. Ohne euch hätte ich es nicht geschafft.

Inhaltsverzeichnis

| | |
|--|-----------|
| 1. Einleitung | 1 |
| 1.1. Motivation | 1 |
| 1.2. Aufbau der Arbeit | 2 |
| 2. Notation, Modellannahmen | 3 |
| 2.1. Notation | 3 |
| 2.2. Modellannahmen | 5 |
| 2.2.1. Modell | 5 |
| 2.3. Hypothesen | 7 |
| 2.3.1. LD-F1 | 7 |
| 2.3.2. LD-F2 | 7 |
| 2.3.3. Allgemeines LD-Fm-Modell | 8 |
| 3. Beispiele | 11 |
| 3.1. α -Amylase Studie | 11 |
| 3.2. Cortisol-Konzentration im Blutplasma | 12 |
| 4. Bekannte Verfahren unter Normalverteilung | 15 |
| 4.1. Hotelling's T^2 -Test | 15 |
| 4.2. Box-Approximation | 16 |
| 4.2.1. ANOVA-Typ Statistik | 17 |
| 4.2.2. Geisser-Greenhouse Statistik | 17 |
| 4.2.3. Schätzer für den Freiheitsgrad f | 18 |
| 4.2.4. ANOVA-Typ Statistik nach Werner | 18 |
| 4.2.5. Geisser-Greenhouse Statistik nach Becker | 19 |
| 4.3. Notwendigkeit dimensionsstabiler Schätzer | 20 |
| 5. ANOVA-Typ Statistik ohne Normalverteilung | 21 |
| 5.1. Grundlagen | 22 |
| 5.2. Ansatz nach Geisser-Greenhouse | 22 |
| 5.3. Box-Approximation eines Quotienten von Zufallsvariablen | 24 |
| 5.3.1. Kovarianz zweier quadratischer Formen | 25 |
| 5.3.2. Varianz einer quadratischen Form | 26 |
| 5.4. Anwendung der Theorie | 27 |

| | |
|---|-----------|
| 5.4.1. Taylor-Approximation | 30 |
| 5.4.2. Box-Approximation | 31 |
| 5.5. Die neue ANOVA-Typ Statistik | 33 |
| 6. Die Schätzer B_1 und B_2 | 35 |
| 6.1. Quadrat- und Bilinearformen | 35 |
| 6.1.1. Darstellungssätze | 37 |
| 6.2. Nachweis der Schätzereigenschaften | 43 |
| 7. ANOVA-Typ Statistik nach Werner ohne Normalverteilung | 45 |
| 7.1. Der Schätzer B_0 | 46 |
| 7.2. Momente der quadratischen Formen | 47 |
| 7.2.1. Taylor-Approximation | 48 |
| 7.2.2. Box-Approximation | 49 |
| 7.3. Die ANOVA-Typ Statistik nach Werner | 51 |
| 8. Simulationen der Statistiken | 53 |
| 8.1. Simulationstechniken | 53 |
| 8.1.1. Niveau und Power | 53 |
| 8.1.2. Kovarianzstrukturen | 55 |
| 8.1.3. Struktur der Zufallszahlen | 55 |
| 8.2. Niveau | 58 |
| 8.3. Power | 62 |
| 9. Software | 65 |
| 9.1. Einbinden und Aufrufen der Makros | 65 |
| 10. Auswertung der Beispiele | 67 |
| 10.1. α -Amylase Studie | 67 |
| 10.2. Cortisol-Konzentration im Blutplasma | 68 |
| 11. Zusammenfassung und Ausblick | 69 |
| A. Anhang | 71 |
| A.1. Makro | 71 |
| A.2. Beweise | 76 |
| A.3. Verwendete Definitionen, Lemmata und Sätze | 88 |
| Literaturverzeichnis | 91 |

Abbildungsverzeichnis

| | |
|---|----|
| 3.1. α -Amylase: links Histogramm, rechts Median-Plots | 11 |
| 3.2. Median-Plots der Cortisol-Konzentration | 12 |
| 8.1. Niveau: Exponentialverteilung | 58 |
| 8.2. Niveau: Normalverteilung | 59 |
| 8.3. Niveau: Log-Normalverteilung | 60 |
| 8.4. Niveau: Gleichverteilung | 60 |
| 8.5. Niveau: Bernulli-Verteilung | 61 |
| 8.6. Ein-Punkt-Power | 62 |
| 8.7. Trend-Power | 62 |

Tabellenverzeichnis

| | |
|--|----|
| 4.1. Simulationen: Normalverteilung zum 95% -Quantil | 20 |
| 10.1. Auswertung der α -Amylase | 67 |
| 10.2. Auswertung der Cortisol-Konzentration | 68 |

1. Einleitung

1.1. Motivation

Die vorliegende Arbeit beschäftigt sich mit der Modellierung und Auswertung von metrischen Datensätzen, in denen verbundene Messungen in Form von Messwiederholungen (*repeated measures*) pro Individuum erhoben werden. Diese Art von Messungen tritt zum Beispiel in klinischen Studien auf, in denen Zeitverläufe (von z.B. Krankheiten oder Wundheilungen) betrachtet werden oder ein Endpunkt unter verschiedenen Bedingungen am selben Individuum gemessen wird. Datensätze dieser Art erhalten zusätzlich die Bezeichnung hochdimensional, wenn die Dimension d (*Anzahl der Messwiederholungen*) den Stichprobenumfang n (*Anzahl der Individuen*) übersteigt. In solchen Datensätzen ist die für viele Teststatistiken benötigte Voraussetzung, dass der Stichprobenumfang größer ist als die Zahl der Messwiederholungen, verletzt.

In der Literatur sind bereits Verfahren für solche Versuchsdesigns bekannt.

Unter Annahme der Normalverteilung der Messwerte konnten Werner (2004) und Becker (2010) Teststatistiken entwickeln, welche für eine beliebige Anzahl an Messwiederholungen ($n < d$ oder $n \geq d$) sehr gute Approximationen liefern. Für asymptotische Resultate ging lediglich n gegen Unendlich und die Anzahl der Messwiederholungen d wurde als beliebig, aber fest angesehen. Ebenfalls unter Normalverteilung gelang Srivastava (2009) die Herleitung einer Teststatistik unter der Annahme, dass n und d gleichzeitig gegen Unendlich gehen.

In der Praxis treten jedoch sehr häufig Datensätze auf, in denen die Daten keiner Normalverteilung folgen. Beispielsweise folgen Reaktionszeiten oder Flächen unter Dosis-Wirkungs-Kurven keiner Normalverteilung, weil diese Daten nicht symmetrisch, sondern schief verteilt sind. Die Annahme der Normalverteilung ist bei solchen Datensätzen also verletzt, sodass die statistische Auswertung und die daraus gezogenen Schlussfolgerungen falsch sein können.

Weitere Verfahren, welche ohne Annahme der Normalverteilung auskommen, benötigen entweder, dass n und d gleichzeitig gegen Unendlich gehen (Bai und Saranadasa, (1996)) oder halten den Stichprobenumfang n fest und lassen nur die Anzahl der Messwiederholungen d gegen Unendlich gehen (Akritas und Wang, (2010)).

Um aber eine Statistik entwickeln zu können, welche schon für eine geringe Anzahl an Messwiederholungen ($d < n$) eine gute Approximation liefert und für eine größere Anzahl ($n \leq d$) nicht schlechter wird, ist eine Herleitung unter der Annahme, dass d gegen Unendlich geht nicht erstrebenswert. Weiterhin ist zu beachten, dass nur in den Verfahren von Werner (2004) und Becker (2010) eine faktorielle Struktur auf den Messwiederholungen zulässig ist.

Das Ziel dieser Arbeit wird es daher sein, die bekannten Verfahren nach Werner (2004) und Becker (2010) dahingehend zu erweitern, dass eine explizite Annahme der Normalverteilung nicht mehr notwendig ist. Dabei sollen analog zu Werner und Becker alle asymptotischen Resultate für beliebiges, aber festes d und für n gegen Unendlich gezeigt werden.

Eine Übersicht der Arbeit wird im folgenden Abschnitt gegeben.

1.2. Aufbau der Arbeit

In Kapitel 2 werden die verwendeten Notationen und Modellannahmen eingeführt. In Kapitel 3 wird anhand von zwei motivierenden Beispielen die Problematik der Arbeit erläutert. In Kapitel 4 werden aus der Literatur bekannte Verfahren unter Normalverteilung vorgestellt und verglichen. In Kapitel 5 wird die neue ANOVA-Typ Statistik ohne Annahme der Normalverteilung hergeleitet, sowie die dafür notwendigen Darstellungssätze für Varianzen und Kovarianzen von quadratischen Formen entwickelt. In Kapitel 6 wird gezeigt, dass die unter Normalverteilung existierenden Schätzer auch ohne Annahme selbiger die gewünschten Eigenschaften wie Erwartungstreue, Konsistenz und Dimensionsstabilität besitzen. In Kapitel 7 folgt dann die Herleitung der ANOVA-Typ Statistik nach Werner ohne Annahme der Normalverteilung. In den Herleitungen der Kapitel 5 bis 7 wird die Anzahl der Messwiederholungen stets als beliebig, aber fest angenommen. Für die asymptotischen Resultate wird lediglich gefordert, dass die Anzahl der Individuen gegen unendlich geht. Danach werden in Kapitel 8 die Simulationsergebnisse zum Niveau und zur Power der beiden hergeleiteten Verfahren im Vergleich zu älteren Statistiken dargestellt. Mit Hilfe der in Kapitel 9 vorgestellten SAS-Makros werden im darauf folgenden Kapitel 10 die zu Anfang erwähnten Beispieldatensätze ausgewertet. Zusammenfassung und Ausblick bilden den Abschluss dieser Arbeit.

2. Notation, Modellannahmen

2.1. Notation

In diesem Abschnitt werden die grundlegenden Notationen eingeführt, die in den folgenden Kapiteln verwendet werden.

Konstante SKALARE werden mit lateinischen Kleinbuchstaben bezeichnet.

VEKTOREN sollen im Allgemeinen mit kleinen, fettgedruckten, lateinischen Buchstaben bezeichnet werden. Des Weiteren sei \mathbf{e}_n der n -dimensionale Einheitsvektor, $\mathbf{1}_n$ der n -dimensionale Einser-Vektor und $\mathbf{0}_n$ der n -dimensionale Null-Vektor.

MATRIZEN werden mit großen fettgedruckten lateinischen Buchstaben bezeichnet und speziell sei \mathbf{I}_n die $n \times n$ Einheitsmatrix, $\mathbf{J}_n = \mathbf{1}_n \cdot \mathbf{1}'_n$ die Einser-Matrix und $\mathbf{P}_n = \mathbf{I}_n - \frac{1}{n} \mathbf{J}_n$ die zentrierende Matrix und ein Projektor.

Sei \mathbf{A}' die TRANSPONIERTE der Matrix \mathbf{A} und bezeichne $Sp(\mathbf{A})$ die SPUR der Matrix \mathbf{A} , sowie $diag(\mathbf{A})$ deren DIAGONALE als Vektor geschrieben. Weiterhin bezeichne $r(\mathbf{A})$ den RANG von \mathbf{A} .

ZUFALLSVARIABLEN werden mit lateinischen Großbuchstaben aus dem Ende des Alphabets bezeichnet und ZUFALLSVEKTOREN werden zusätzlich in Fettdruck dargestellt (Z_{kl}, \mathbf{Z}_k).

Alle INDIZES von Vektoren, Matrizen und Zufallsvariablen sollen mit kleinen lateinischen Buchstaben wie zum Beispiel k, l, r, s bezeichnet werden.

Darüber hinaus werden noch einige Matrizenoperationen angeführt, die in der gesamten Arbeit von essentieller Bedeutung sind.

Die KRONECKER-SUMME $\mathbf{A} \oplus \mathbf{B} = \left(\begin{array}{c|c} \mathbf{A} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{B} \end{array} \right)$ zweier Matrizen \mathbf{A} und \mathbf{B} wird mit \oplus beschrieben.

Das KRONECKER-PRODUKT $\mathbf{A} \otimes \mathbf{B} = (a_{ij} \mathbf{B})_{ij} = \begin{pmatrix} a_{11} \mathbf{B} & \cdots & a_{1n} \mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{n1} \mathbf{B} & \cdots & a_{nn} \mathbf{B} \end{pmatrix}$ zweier Matrizen \mathbf{A} und \mathbf{B} wird mit \otimes beschrieben, wobei a_{ij} die Elemente der Matrix \mathbf{A} sind.

KONVERGENZEN

Für eine Folge von Zufallsvariablen (X_n) bezeichne $X_n \xrightarrow{p} X$ die stochastische Konvergenz, $X_n \xrightarrow{\mathcal{L}} X$ die Verteilungskonvergenz und $X_n \xrightarrow{\mathcal{L}_2} X$ die \mathcal{L}_2 -Konsistenz der Zufallsvariablen X_n gegen X . Sei (Y_n) ebenfalls eine Folge von Zufallsvariablen, dann bezeichne $X_n \doteq Y_n$ die asymptotische Äquivalenz (siehe Definition A.3.5).

2.2. Modellannahmen

In vielen Arbeiten, die sich mit hochdimensionalen Repeated-Measures-Versuchsplänen beschäftigten, wurde die Annahme der Normalverteilung der Messgrößen zugrunde gelegt (siehe z.B. Werner (2004), Srivastava (2009) oder Becker (2010)). Diese Art von Modellen hatte folgende Gestalt:

Seien

$$\mathbf{X}_k = (X_{k1}, \dots, X_{kd})' \sim N(\boldsymbol{\mu}, \mathbf{V}), \quad k = 1, \dots, n,$$

unabhängig multivariat normalverteilte Zufallsvektoren mit Erwartungswertvektor $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)'$ und Kovarianzmatrix \mathbf{V} .

In dieser Arbeit soll auf die explizite Annahme der Normalverteilung der Zufallsvektoren verzichtet werden. Dabei ist zu beachten, dass auch ohne die Annahme der Normalverteilung ein parametrisches Modell zugrunde gelegt wird und die Hypothesen weiterhin über den Erwartungsvektor $\boldsymbol{\mu}$ definiert werden. Die folgenden Betrachtungen beschränken sich daher auf Verteilungen von metrischen Daten (keine Scores).

Die unter diesen Annahmen hergeleiteten Teststatistiken sind auf eine breitere Klasse von Verteilungen anwendbar und beschränken sich nicht mehr auf die Normalverteilung.

Der Ansatz für das allgemeine Modell ist ähnlich dem im Paper von Bai und Saranadasa (1996) S. 311-329 gewählten Modell, wobei nur ein Teil der dort getroffenen Annahmen benötigt wird. Dieser Ansatz wird im nächsten Abschnitt diskutiert.

2.2.1. Modell

Es werden

$$\mathbf{X}_k = \boldsymbol{\Gamma} \mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}, \quad k = 1, \dots, n, \quad (2.1)$$

unabhängig identisch verteilte Zufallsvektoren betrachtet, mit $\boldsymbol{\Gamma} \in \mathbb{R}^{d \times d}$ beliebig, aber pro Versuchsdesign fest und $\boldsymbol{\Gamma} \boldsymbol{\Gamma}' = \mathbf{S}$.

Die $\mathbf{Z}_k = (Z_{k1}, \dots, Z_{kd})'$, $k = 1, \dots, n$, seien ebenfalls unabhängig identisch verteilte Zufallsvektoren, deren Komponenten Z_{ks} unabhängig sind, für $s = 1, \dots, d$. Zusätzlich seien $E(\mathbf{Z}_k) = \mathbf{0}$, $Var(\mathbf{Z}_k) = \mathbf{I}_d$ und die vierten Momente der einzelnen Komponenten wie folgt beschränkt: $E(Z_{ks}^4) \leq \gamma < \infty \forall k = 1, \dots, n, s = 1, \dots, d$.

Dann bezeichne $\boldsymbol{\Gamma} \mathbf{Z}_k$ den Versuchsfehler und $\mathbf{S} = Cov(\boldsymbol{\Gamma} \mathbf{Z}_1)$ die zugehörige Kovarianzmatrix sowie $E_k \mathbf{1}_d$ die interindividuelle Streuung für das k -te Subjekt, wobei die Zufallsvariablen der interindividuellen Streuung E_k ebenfalls unabhängig identisch verteilt sind für $k = 1, \dots, n$, sowie nach Konstruktion unabhängig von Z_{ks}

$\forall k = 1, \dots, n, s = 1, \dots, d$, mit $E(E_k) = 0$ und $Var(E_k) = E(E_k^2) = \sigma_{E_k}^2 < \infty$.

Aus dieser Konstruktion ergibt sich $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)' = E(\mathbf{X}_1)$ als Erwartungswertvektor und $\mathbf{V} = \mathbf{S} + \sigma_{E_k}^2 \mathbf{J}_d = Cov(\mathbf{X}_1)$ als Kovarianzmatrix der Zufallsvektoren \mathbf{X}_k , $k = 1, \dots, n$, wobei auch $d > n$ zulässig ist.

Die Struktur der Kovarianzmatrix des Versuchsfehlers \mathbf{S} hängt von der Art des Versuchsaufbaus ab. Es werden verschiedene Arten von Kovarianz-Strukturen unterschieden, die verschiedene Arten von Abhängigkeiten darstellen, wobei für die in dieser Arbeit betrachteten Repeated-Measures-Designs am häufigsten Compound Symmetry und Autoregression auftreten. Diese werden im Folgenden diskutiert.

Definition 2.2.1 (Compound Symmetry (CS))

Sei $\mathbf{X} = (X_1, \dots, X_d)'$ ein Vektor von Zufallsvariablen. Dann heißt die Kovarianzstruktur compound symmetric, wenn gilt:

$$\begin{aligned} Var(X_i) &= \sigma^2 \quad \forall i = 1, \dots, d, \\ Cov(X_i, X_j) &= \tau \quad \forall i \neq j = 1, \dots, d. \end{aligned}$$

Die Abhängigkeiten zwischen den Zufallsvariablen sind immer gleich.

Definition 2.2.2 (Autoregression (AR))

Sei $\mathbf{X} = (X_1, \dots, X_d)'$ ein Vektor von Zufallsvariablen, dann heißt die Kovarianzstruktur autoregressiv, wenn gilt:

$$Cov(X_i, X_j) = \rho^{|i-j|} \cdot \sigma^2 \quad \forall i, j = 1, \dots, d, \rho \in (0,1).$$

Je weiter zwei Zufallsvariablen auseinander liegen, desto geringer ist ihre Abhängigkeit zueinander.

Weiterhin kann eine beliebige faktorielle Struktur auf den wiederholten Messungen betrachtet werden. Beim LD-F1-Aufbau der Daten gibt es d Messwiederholungen (*repeated measures*) pro Individuum in einer homogenen Gruppe von n unabhängigen Individuen, dadurch bekommen die Zufallsvektoren eine einfache Indizierung:

$$\mathbf{X}_k = (X_{k1}, \dots, X_{kd})', \quad k = 1, \dots, n.$$

Beim LD-F2-Aufbau gibt es b wiederholte Messungen pro Individuum in einer Gruppe von n unabhängigen Individuen, die zusätzlich unter a verschiedenen Bedingungen durchgeführt werden. Die Zufallsvektoren erhalten somit eine zweifache Strukturierung:

$$\mathbf{X}_k = (X_{k11}, \dots, X_{k1b}, \dots, X_{kab})', \quad k = 1, \dots, n.$$

Wobei zu beachten ist, dass insgesamt $d = a \cdot b$ Messungen desselben Endpunktes an jedem Individuum vorgenommen werden.

Allgemein können Versuche mit beliebig vielen Faktoren F_1, \dots, F_m auf den Messwiederholungen untersucht werden, der Vektor $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)'$ erhält dann eine zusätzliche m -fache Strukturierung. Dabei muss stets zwischen wiederholten Messungen und multiplen Endpunkten unterschieden werden. Das Hauptaugenmerk dieser Arbeit liegt allerdings auf den wiederholten Messungen (*repeated measures*).

2.3. Hypothesen

2.3.1. LD-F1

In einem LD-F1-Versuchsplan ist zunächst die globale Hypothese von Interesse, nämlich die Frage, ob ein globaler Zeiteffekt vorliegt. Mathematisch formuliert lautet die Hypothese "kein Zeiteffekt (T)" wie folgt:

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_d \Leftrightarrow \mu_i = \bar{\mu}, \quad i = 1, \dots, d.$$

In Matrizenschreibweise lässt sich die Hypothese folgendermaßen darstellen:

$$H_0(T) : \left(\mathbf{I}_d - \frac{1}{d} \mathbf{J}_d \right) \boldsymbol{\mu} = \mathbf{P}_d \boldsymbol{\mu} \quad \text{mit } \boldsymbol{\mu} = (\mu_1, \dots, \mu_d)'.$$

2.3.2. LD-F2

Im LD-F2-Aufbau sind drei Hypothesen interessant: "kein Zeiteffekt (T)", "kein Behandlungs- oder Bedingungseffekt (B)" und "keine Wechselwirkung zwischen Zeit und Behandlung bzw. Bedingung (TB)". Mathematisch und in Matrizenschreibweise lauten die Hypothesen dann:

$$H_0(T) : \bar{\mu}_{\cdot 1} = \dots = \bar{\mu}_{\cdot b} \Leftrightarrow \bar{\mu}_{\cdot i} = \bar{\mu}_{\cdot\cdot}, \quad i = 1, \dots, b,$$

$$\left(\frac{1}{a} \mathbf{1}'_a \otimes \mathbf{P}_b \right) \boldsymbol{\mu} = \mathbf{H}_T \boldsymbol{\mu} = \mathbf{0}$$

$$H_0(B) : \bar{\mu}_{1\cdot} = \dots = \bar{\mu}_{a\cdot} \Leftrightarrow \bar{\mu}_{j\cdot} = \bar{\mu}_{\cdot\cdot}, \quad j = 1, \dots, a,$$

$$\left(\mathbf{P}_a \otimes \frac{1}{b} \mathbf{1}'_b \right) \boldsymbol{\mu} = \mathbf{H}_B \boldsymbol{\mu} = \mathbf{0}$$

$$H_0(TB) : \mu_{ji} + \bar{\mu}_{\cdot\cdot} = \bar{\mu}_{j\cdot} + \bar{\mu}_{\cdot i}, \quad j = 1, \dots, a, \quad i = 1, \dots, b$$

$$(\mathbf{P}_a \otimes \mathbf{P}_b) \boldsymbol{\mu} = \mathbf{H}_{TB} \boldsymbol{\mu} = \mathbf{0}.$$

2.3.3. Allgemeines LD-Fm-Modell

Für die Herleitung der Theorie wird im Folgenden immer die kanonische Formulierung einer Hypothesenmatrix \mathbf{H} mittels eines Projektors \mathbf{T} verwendet. Dieser ist idempotent und symmetrisch. Aus jeder Hypothesenmatrix \mathbf{H} kann dieser Projektor \mathbf{T} über die Beziehung $\mathbf{T} = \mathbf{H}'(\mathbf{H}\mathbf{H}')^{-1}\mathbf{H}$ konstruiert werden. Denn unter H_0 gilt für jede Hypothesenmatrix \mathbf{H} : $\mathbf{T}\boldsymbol{\mu} = \mathbf{0} \Leftrightarrow \mathbf{H}\boldsymbol{\mu} = \mathbf{0}$. Es genügt somit, alle Teststatistiken und Verfahren für Hypothesen der Form $\mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ herzuleiten und auf gewünschte Modelle anzuwenden, indem der entsprechende Projektor $\mathbf{T} = \mathbf{H}'(\mathbf{H}\mathbf{H}')^{-1}\mathbf{H}$ mittels der zugehörigen Hypothesenmatrix \mathbf{H} gebildet wird. \mathbf{T} ist dann sogar die kanonische Wahl einer Hypothesenmatrix (siehe Boysen (2002)). Weiterhin sei angemerkt, dass in allen folgenden Betrachtungen die Hypothesenmatrix \mathbf{H} als eine Kontrastmatrix definiert ist. Diese Eigenschaft überträgt sich auf den Projektor \mathbf{T} und speziell folgt daraus: $\mathbf{T} \cdot \mathbf{1}_d = \mathbf{0}$.

Die Herleitung der Theorie hängt dann nicht von der Struktur des Modells ab, solange der Projektor \mathbf{T} aus der Hypothesenmatrix \mathbf{H} gebildet werden kann. Daher wird zur übersichtlicheren Darstellung der Theorie das LD-F1 Modell verwendet, sodass die Zufallsvektoren \mathbf{X}_k in den folgenden Betrachtungen nur eine einfache Indizierung besitzen.

Weiterhin werden im Folgenden häufig die, auf den Hypothesenraum projizierten, Zufallsvektoren \mathbf{TX}_k benötigt, welche mit:

$$\mathbf{Y}_k = \mathbf{TX}_k, \quad k = 1, \dots, n \quad (2.2)$$

bezeichnet werden sollen.

Die Eigenschaften der Zufallsvektoren \mathbf{Y}_k , unter den Modellvoraussetzungen aus (2.1), werden in der folgenden Proposition vorgestellt.

Proposition 2.3.1 *Die Zufallsvektoren $\mathbf{X}_k = (X_{k1}, \dots, X_{kd})'$, $k = 1, \dots, n$, seien unabhängig identisch verteilt und erfüllen die Modellvoraussetzungen aus (2.1) mit $\mathbf{X}_k = \mathbf{\Gamma Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}$. Sei weiterhin $\mathbf{T} = \mathbf{H}'(\mathbf{H H}')^{-1} \mathbf{H}$ der in Abschnitt 2.3.3 definierte Projektor und die zugehörige Hypothesenmatrix \mathbf{H} eine Kontrastmatrix. Dann lassen sich die \mathbf{Y}_k wie folgt darstellen:*

$$\begin{aligned} \mathbf{Y}_k &= \mathbf{TX}_k = \mathbf{T \Gamma Z}_k + E_k \mathbf{T 1}_d + \mathbf{T \boldsymbol{\mu}} = \mathbf{T \Gamma Z}_k + \mathbf{T \boldsymbol{\mu}} \\ \mathbf{Y}_k &= (Y_{k1}, Y_{k2}, \dots, Y_{kd})', \text{ mit} \\ \text{Cov}(\mathbf{Y}_k) &= \mathbf{T V T} = \mathbf{T} (\mathbf{S} + \sigma_{E_k}^2 \mathbf{J}_d) \mathbf{T} = \mathbf{T S T} = \boldsymbol{\Sigma}, \quad k = 1, \dots, n. \end{aligned}$$

Unter Hypothese $H_0 : \mathbf{T \boldsymbol{\mu}} = \mathbf{0}$ folgt dann:

$$\begin{aligned} \mathbf{Y}_k &= \mathbf{TX}_k = \mathbf{T \Gamma Z}_k \\ E_{H_0}(\mathbf{Y}_k) &= \mathbf{0}, \quad k = 1, \dots, n. \end{aligned}$$

Im nächsten Kapitel werden die vorgestellten Modelle anhand von zwei motivierenden Beispielen weiter diskutiert.

3. Beispiele

3.1. α -Amylase Studie

In dieser Studie aus der HNO-Klinik (siehe Brunner (2002) S.9) wurde die Konzentration der α -Amylase in der Speichelflüssigkeit von 14 Patienten untersucht. Die Speichelproben wurden an 2 Tagen (*Montag, Donnerstag*) zu jeweils 4 Zeitpunkten (8 a.m., 12 a.m., 5 p.m., 9 p.m.) bei allen Probanden entnommen. Die Fragestellung bestand darin, ob sich die α -Amylase im Tagesverlauf verändert, wobei zusätzlich ein Unterschied im α -Amylase Profil direkt nach dem Wochenende gegenüber der Mitte der Woche vermutet wurde. Es liegt also ein zweifaktorielles Repeated-Measures-Design vor, welches mit 14 Probanden und 8 Messwiederholungen je Proband nicht hochdimensional ist. Die Verteilung der α -Amylase wird im Histogramm in Abbildung 3.1 dargestellt, wobei hier die faktorielle Struktur vernachlässigt wurde. In Abbildung 3.1 (rechts) sind Median-Plots der α -Amylase dargestellt.

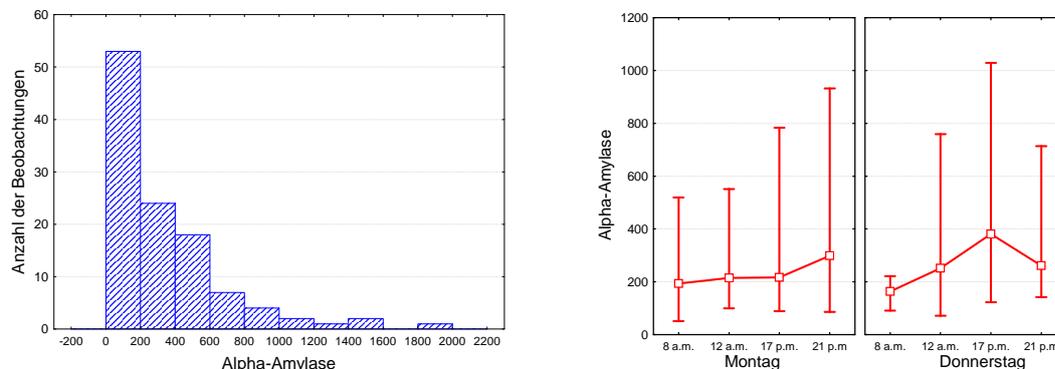


Abbildung 3.1.: α -Amylase: links Histogramm, rechts Median-Plots

Aus dem Histogramm (Abbildung 3.1) wird ersichtlich, dass eine Annahme der Normalverteilung für die stetigen α -Amylase Daten nicht sinnvoll ist. Weiterhin zeigen die Median-Plots (Abbildung 3.1), dass der vermutete Effekt zwischen den Tagen gut zu beobachten ist, wohingegen eine Veränderung im Tagesverlauf nur am Donnerstag zu erkennen ist.

3.2. Cortisol-Konzentration im Blutplasma

In diesem Versuch wurde die Cortisol-Konzentration im Blutplasma von 12 Marathonläufern untersucht (siehe Brunner (2002) S. 10), die eine abrupte Trainingspause über zwei Wochen einhalten mussten. Im Falle solch einer radikalen Trainingsunterbrechung treten sogenannte „Sportentzugerscheinungen“ auf, die mit Unbehagen, Schlaflosigkeit und Verdauungsstörungen einhergehen können. Ziel dieses Versuches war es, die Veränderungen der Cortisol-Konzentration bei „Sportentzugerscheinungen“ zu untersuchen. Weiterhin war die Frage von Interesse, ob die Verabreichung von *m-CCP* die erwarteten Schwankungen in der Cortisol-Konzentration verringern und somit gegebenenfalls die Symptome des Sportentzugs reduzieren könnte. Dazu wurde allen 12 Probanden vor und nach der zweiwöchigen Trainingspause (*vor*, *nach*) die Substanz *m-CCP* verabreicht und anschließend wurden jeweils sieben Messungen (*0 min.*, *30 min.*, *60 min.*, *90 min.*, *120 min.*, *180 min.*, *240 min.*) der Cortisol-Konzentration in kurzen Abständen durchgeführt. Nach einer Auswaschzeit vor und nach der Trainingspause wurde die Cortisol-Konzentration, ohne Behandlung (*Placebo*), noch einmal in den gleichen Abständen sieben Mal gemessen. Damit ergibt sich ein 3-faktorielles hochdimensionales Repeated-Measures-Design, da an 12 Probanden 28-mal die Cortisol-Konzentration gemessen wurde. In Abbildung 3.2 sind Median-Plots der Cortisol-Konzentration dargestellt.

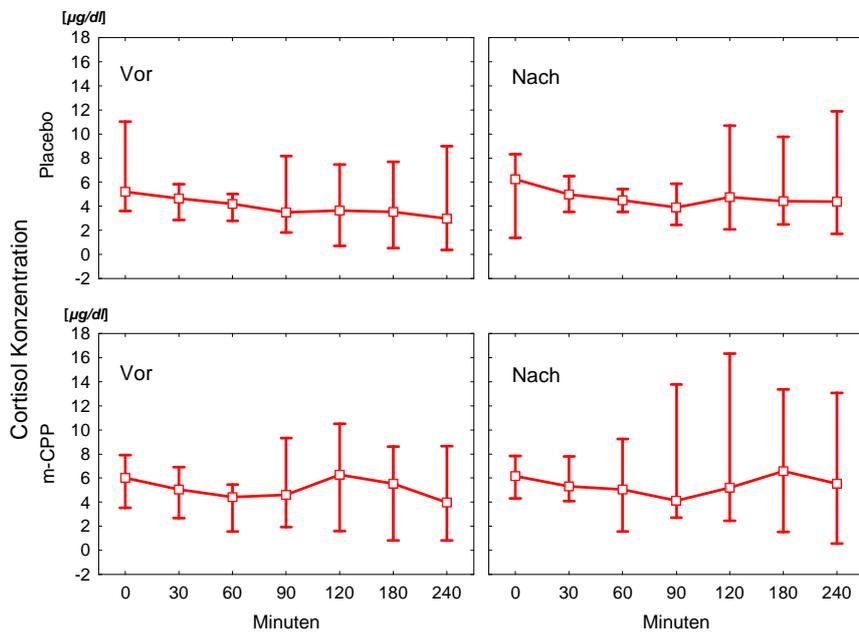


Abbildung 3.2.: Median-Plots der Cortisol-Konzentration

Aus den Median-Plots (Abbildung 3.2) lässt sich die vermutete Veränderung der Cortisol-Konzentration durch die Verabreichung von m-CCP ablesen, sowie eine leichte Veränderung vor und nach der Trainingspause unter beiden Behandlungen erkennen. Eine deutliche Veränderung in den wiederholten Messungen an den jeweiligen Tagen kann aber nicht beobachtet werden. Die Daten sind deutlich schief verteilt, sodass auch hier eine Annahme der Normalverteilung nicht sinnvoll ist.

Beide Beispiele werden in Kapitel 11 ausgewertet.

3. Beispiele

4. Bekannte Verfahren unter Normalverteilung

In diesem Kapitel werden einige grundlegende Teststatistiken vorgestellt, sowie ihre Weiterentwicklungen auf hochdimensionale Versuchspläne erläutert.

Seien dafür

$$\mathbf{X}_k = (X_{k1}, \dots, X_{kd})' \sim N(\boldsymbol{\mu}, \mathbf{V}), \quad k = 1, \dots, n,$$

unabhängig identisch verteilte Zufallsvektoren mit entsprechendem Erwartungswertvektor $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)'$ und Kovarianzmatrix \mathbf{V} . Weiterhin sei $\bar{\mathbf{X}} = (\bar{X}_{\cdot 1}, \dots, \bar{X}_{\cdot d})'$ der Vektor der arithmetischen Mittel und $\hat{\mathbf{V}}_n$ die empirische Kovarianzmatrix. Die betrachteten Hypothesen seien von der Form $H_0 : \mathbf{H}\boldsymbol{\mu} = \mathbf{0}$ mit $r(\mathbf{H}) = p \leq d$, wobei \mathbf{H} und der entsprechende Projektor $\mathbf{T} = \mathbf{H}'(\mathbf{H}\mathbf{H}')^{-1}\mathbf{H}$ Kontrastmatrizen sind.

4.1. Hotelling's T^2 -Test

Eine der klassischen Teststatistiken, welche auch in Repeated-Measures-Designs ihre Anwendung findet, ist der von Hotelling (1931) hergeleitete T^2 -Test. Hierfür wird der Vektor $\mathbf{H}\bar{\mathbf{X}}$ betrachtet und von beiden Seiten an die Inverse der empirischen Kovarianzmatrix $(\mathbf{H}\hat{\mathbf{V}}_n\mathbf{H}')^{-1}$ multipliziert:

$$t^2 = n \cdot (\mathbf{H}\bar{\mathbf{X}})' (\mathbf{H}\hat{\mathbf{V}}_n\mathbf{H}')^{-1} (\mathbf{H}\bar{\mathbf{X}}) \sim T^2(p, n-1, \delta). \quad (4.1)$$

Sei $T^2(p, n-1, \delta)$ die Hotelling's T^2 -Verteilung mit dem Nichtzentralitätsparameter $\delta = n \cdot (\mathbf{H}\boldsymbol{\mu})' (\mathbf{H}\mathbf{V}\mathbf{H}')^{-1} (\mathbf{H}\boldsymbol{\mu})$, welcher unter $H_0 : \mathbf{H}\boldsymbol{\mu} = \mathbf{0}$ verschwindet.

Solange die Anzahl der Messwiederholungen d kleiner als der Stichprobenumfang n ist, kann diese Statistik berechnet werden und liefert eine gute Approximation für die Verteilung von t^2 durch eine T^2 -Verteilung. Liegt hingegen ein hochdimensionales Versuchsdesign ($d > n$) vor, dann ist die empirische Kovarianzmatrix $\hat{\mathbf{V}}_n$ singulär und es ist nicht möglich deren Inverse $\hat{\mathbf{V}}_n^{-1}$ zu berechnen. Das bedeutet für die Betrachtung von speziell hochdimensionalen Versuchsplänen kann der Hotellings T^2 -Test

nicht berechnet werden und wird daher im weiteren Verlauf dieser Arbeit nicht mehr betrachtet.

Anders verhält es sich mit den im Nachfolgenden vorgestellten Verfahren, welche in hoch- und niedrigdimensionalen Versuchsdesigns bestimmt werden können.

4.2. Box-Approximation

Das Ziel ist es, Statistiken für das Testen von Hypothesen der Form $\mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ herzuleiten. Dazu wird zuerst die quadratische Form der Mittelwertsvektoren $Q_n^* = n \cdot \bar{\mathbf{X}}' \mathbf{T} \bar{\mathbf{X}}$ betrachtet und versucht deren Verteilung zu bestimmen.

Dies gelingt approximativ, da sich quadratische Formen in normalverteilten Zufallsvariablen als eine gewichtete Summe von χ_1^2 -Verteilungen darstellen lassen. Diese Verteilung der Linearkombinationen $\sum_{i=1}^d \lambda_i C_i$ kann dann durch die Verteilung einer mit g gestreckten χ_f^2 -Verteilung approximiert werden, so dass die ersten beiden Momente von $\sum_{i=1}^d \lambda_i C_i$ und $g \cdot U \sim g\chi_f^2$ übereinstimmen.

Dieses Verfahren, eine Verteilung zu approximieren, wird analog zu vorangegangenen Arbeiten (Werner, (2004), Becker, (2010)) „Box-Approximation“ genannt. Für eine explizite Herleitung siehe Box I (1954) und Box II (1954). Dann gelten folgende Beziehungen:

$$E \left(\sum_{i=1}^d \lambda_i C_i \right) = \sum_{i=1}^d \lambda_i \stackrel{!}{=} E(g \cdot U) = g \cdot f = Sp(\boldsymbol{\Sigma})$$

$$Var \left(\sum_{i=1}^d \lambda_i C_i \right) = 2 \cdot \sum_{i=1}^d \lambda_i^2 \stackrel{!}{=} Var(g \cdot U) = 2 \cdot g^2 \cdot f = 2 \cdot Sp(\boldsymbol{\Sigma}^2).$$

Es folgt daraus $f = [Sp(\boldsymbol{\Sigma})]^2 / Sp(\boldsymbol{\Sigma}^2)$ und $g = Sp(\boldsymbol{\Sigma}^2) / Sp(\boldsymbol{\Sigma})$. Somit hat die quadratische Form Q_n^* die approximative Verteilung:

$$\frac{n \cdot \bar{\mathbf{X}}' \mathbf{T} \bar{\mathbf{X}}}{Sp(\boldsymbol{\Sigma})} = \frac{Q_n^*}{Sp(\boldsymbol{\Sigma})} \dot{\sim} \chi_f^2 / f = F(f, \infty). \quad (4.2)$$

4.2.1. ANOVA-Typ Statistik

Wird in (4.2) zum Schätzen der unbekanntes Kovarianzmatrix die empirische Kovarianzmatrix

$$\widehat{\Sigma}_n = \frac{1}{(n-1)} \sum_{k=1}^n (\mathbf{T}\mathbf{X}_k - \mathbf{T}\bar{\mathbf{X}}.)' (\mathbf{T}\mathbf{X}_k - \mathbf{T}\bar{\mathbf{X}}.)$$

verwendet, so ergibt sich die klassische ANOVA-Typ Statistik (ATS) unter Annahme der Normalverteilung, zum Testen von $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$:

$$A_n^{ATS} = \frac{Q_n^*}{Sp(\widehat{\Sigma}_n)} \dot{\sim} \chi_{\widehat{f}}^2 / \widehat{f}, \quad \widehat{f} = \frac{[Sp(\widehat{\Sigma}_n)]^2}{Sp(\widehat{\Sigma}_n^2)}. \quad (4.3)$$

Wird die empirische Kovarianzmatrix $\widehat{\Sigma}_n$ ebenfalls zum Schätzen von f verwendet, dann ist diese Schätzung nicht erwartungstreu, sondern verzerrt. Die Verzerrung wächst mit steigender Dimension d . Der Grund dafür ist der Schätzer $Sp(\widehat{\Sigma}_n^2)$, welcher mit wachsender Dimension d positiv verzerrt wird. Dieses Problem wurde bereits von Werner (2004) bearbeitet.

4.2.2. Geisser-Greenhouse Statistik

Wird in (4.2) die unbekanntes Spur der Kovarianzmatrix als quadratische Form dargestellt und darauf ebenfalls die Box-Approximation angewendet, so ergibt sich die Statistik von Geisser-Greenhouse (1958) (GG):

$$A_n^{GG} = \frac{Q_n^*}{Sp(\widehat{\Sigma}_n)} \dot{\sim} F(\widehat{f}, (n-1) \cdot \widehat{f}), \quad \widehat{f} = \frac{[Sp(\widehat{\Sigma}_n)]^2}{Sp(\widehat{\Sigma}_n^2)}. \quad (4.4)$$

Auch bei dieser Teststatistik besteht das Problem in der adäquaten Schätzung des Freiheitsgrades $f = [Sp(\boldsymbol{\Sigma})]^2 / Sp(\boldsymbol{\Sigma}^2)$. Wird dafür, analog zur ANOVA-Typ-Statistik (4.3), die empirische Kovarianzmatrix $\widehat{\Sigma}_n$ zur Schätzung von $[Sp(\boldsymbol{\Sigma})]^2$ und $Sp(\boldsymbol{\Sigma}^2)$ verwendet, so ist der Schätzer \widehat{f} verzerrt und die Statistik wird mit wachsender Dimension immer konservativer.

Dieses Problem wurde speziell von Becker (2010) bearbeitet. Daher werden die Ergebnisse von Werner (2004) und Becker (2010) in den nächsten Abschnitten kurz vorgestellt.

Es sei noch angemerkt, dass Geisser und Greenhouse (1958) aus diesem Grund die untere Grenze $f = 1$ als Freiheitsgrad vorgeschlagen haben, welche aber oft zu noch

konservativeren Ergebnissen führt als die Verwendung der empirischen Kovarianzmatrix.

4.2.3. Schätzer für den Freiheitsgrad f

Um die Verzerrung des Freiheitsgradschätzers \hat{f} zu verringern, wurden erwartungstreue, konsistente und dimensionsstabile Schätzer (siehe Definition A.3.4 Dimensionsstabilität) für die Spuren $[Sp(\boldsymbol{\Sigma})]^2$ und $Sp(\boldsymbol{\Sigma}^2)$ der Kovarianzmatrix benötigt.

Die von Werner (2004) hergeleiteten Schätzer B_1 und B_2 sind erwartungstreu und erfüllen unter Annahme der Normalverteilung die ebenfalls geforderten Bedingungen. Weiterhin wurde in Werner (2004) ein erwartungstreuer Schätzer für $Sp(\boldsymbol{\Sigma})$ entwickelt, der im Folgenden mit B_0 bezeichnet werden soll.

Sie sind die arithmetischen Mittel von Quadrat- und Bilinearformen, die wie folgt definiert sind:

$$\begin{aligned} A_k &= \mathbf{X}'_k \mathbf{T} \mathbf{X}_k \\ A_{kl} &= \mathbf{X}'_k \mathbf{T} \mathbf{X}_l, \quad k \neq l. \end{aligned}$$

Mit deren Hilfe ergeben sich die Spurschätzer B_0 , B_1 und B_2 mit folgenden Erwartungswerten:

$$\begin{aligned} B_0 &= \frac{1}{n} \cdot \sum_{k=1}^n A_k & E_{H_0}(B_0) &= Sp(\boldsymbol{\Sigma}) \\ B_1 &= \frac{1}{n(n-1)} \cdot \sum_{k \neq l} A_k \cdot A_l & E_{H_0}(B_1) &= [Sp(\boldsymbol{\Sigma})]^2 \\ B_2 &= \frac{1}{n(n-1)} \cdot \sum_{k \neq l} A_{kl}^2 & E_{H_0}(B_2) &= Sp(\boldsymbol{\Sigma}^2). \end{aligned}$$

4.2.4. ANOVA-Typ Statistik nach Werner

Durch das Einsetzen der zuvor definierten Schätzer in (4.2) ergibt sich die von Werner (2004) (siehe auch Ahmad et al. (2008)) hergeleitete ANOVA-Typ Statistik (ATS-Werner):

$$A_n^{ATS-W} = \frac{Q_n^*}{B_0} \rightsquigarrow \chi_{\hat{f}}^2 / \hat{f}, \quad \hat{f} = \frac{n}{(n-1)} \cdot \frac{B_1}{B_2}. \quad (4.5)$$

4.2.5. Geisser-Greenhouse Statistik nach Becker

Analog werden in der Geisser-Greenhouse Statistik (4.4) die eben definierten Schätzer B_1 und B_2 verwendet, um den Freiheitsgrad f zu schätzen. Daraus ergibt sich die in Becker (2010) hergeleitete Statistik (siehe auch Bunner (2009))

$$A_n^{GG-B} = \frac{Q_n^*}{Sp(\widehat{\Sigma}_n)} \dot{\sim} F(\tilde{f}, (n-1)\tilde{f}), \quad \tilde{f} = \frac{B_1}{B_2 \left(1 + \frac{1}{4n(n-1)}\right)}, \quad (4.6)$$

welche im Folgenden mit GG-Becker Statistik bezeichnet werden soll.

4.3. Notwendigkeit dimensionsstabiler Schätzer

Anhand von Simulationen soll die Notwendigkeit von erwartungstreuen, konsistenten und speziell dimensionsstabilen Schätzern im Repeated-Measures-Design verdeutlicht werden. Wie aus Tabelle 4.1 ersichtlich ist, degenerieren die klassischen Verfahren (ATS und GG) mit wachsender Dimension schon unter Normalverteilung mit Compound Symmetry als Kovarianzstruktur (siehe Definition 2.2.1). Daher kann im Repeated-Measures-Design ohne Annahme der Normalverteilung nicht auf Schätzer verzichtet werden, die neben Erwartungstreue und Konsistenz auch die Eigenschaft der Dimensionsstabilität besitzen.

In den folgenden Kapiteln werden Statistiken mit den geforderten Eigenschaften hergeleitet.

Tabelle 4.1.: Niveau: 95%-Quantile der Statistiken mit $n = 10$ und CS(2,1)

| d | ATS | ATS-Werner | GG | GG-Becker |
|-----|--------|------------|--------|-----------|
| 3 | 0,9318 | 0,9498 | 0,9578 | 0,9457 |
| 5 | 0,9532 | 0,9513 | 0,9725 | 0,9484 |
| 8 | 0,9699 | 0,9519 | 0,9817 | 0,9497 |
| 10 | 0,9747 | 0,9483 | 0,9871 | 0,9464 |
| 20 | 0,9934 | 0,9559 | 0,9969 | 0,9542 |
| 30 | 0,9981 | 0,9473 | 0,9991 | 0,9461 |
| 50 | 0,9997 | 0,9527 | 1,0000 | 0,9514 |
| 100 | 1,0000 | 0,9521 | 1,0000 | 0,9513 |
| 150 | 1,0000 | 0,9501 | 1,0000 | 0,9494 |

Es ist zu beachten, dass alle bisher vorgestellten Statistiken die Normalverteilung der Zufallsvektoren voraussetzen, wobei die Statistik nach ATS-Werner (4.5) diese Voraussetzung nicht zur Herleitung der Teststatistik, sondern lediglich zum Nachweis der genannten Eigenschaften der Schätzer B_1 und B_2 benötigt.

5. ANOVA-Typ Statistik ohne Normalverteilung

Das Ziel dieses Kapitels wird es daher sein, über die beiden vorgestellten Ansätze der ANOVA-Typ Statistik (4.3) und der Geisser-Greenhouse Statistik (4.4) mit Hilfe der in Werner (2004) entwickelten Schätzer eine Teststatistik herzuleiten, die robust unter der Annahme der Normalverteilung ist. Dazu werden die beiden Ansätze im Vorfeld verglichen, um festzustellen, welcher sich besser auf die neuen Modellannahmen erweitern lässt.

Es ist daher zu beachten, dass der Schätzer der Spur der Kovarianzmatrix B_0 , welcher in der Werner-Statistik benötigt wird, nur im Ein-Stichprobenfall ohne Einschränkungen unverzerrt bestimmt werden kann. Dies ist im Mehr-Stichprobenfall nur unter starken Restriktionen hinsichtlich des Versuchsdesigns möglich. So müssten zum Beispiel im Zwei-Stichprobenfall entweder die Kovarianzmatrizen oder die Stichprobenumfänge in beiden Gruppen gleich sein (siehe auch Ahmad (2008)). Es sollte aber eine Statistik hergeleitet werden, die eine Erweiterung auf den Mehr-Stichprobenfall zulässt, ohne eben genannte Modelleinschränkungen zu benötigen.

Dies wurde unter Annahme der Normalverteilung für den Zwei-Stichprobenfall über den Ansatz von Box und Geisser-Greenhouse bereits von Becker (2010) hergeleitet. Daher wird der ATS-Werner Ansatz zwar weiter verfolgt und in einem separaten Kapitel dargestellt, doch das Hauptaugenmerk liegt auf dem Ansatz von Geisser und Greenhouse über die Spur der empirischen Kovarianzmatrix, welcher von jetzt an Gegenstand der weiteren Betrachtung sein wird. Damit werden sich für den in dieser Arbeit betrachteten Ein-Stichprobenfall zwei Teststatistiken ergeben, die abschließend miteinander verglichen werden.

In den nächsten Abschnitten wird die Herleitung einer Teststatistik ohne Annahme der Normalverteilung erläutert.

5.1. Grundlagen

Aufgrund der Tatsache, dass keine explizite Normalverteilung der betrachteten Zufallsvektoren gefordert wurde, können einige der nachfolgenden Ergebnisse nur asymptotisch mit Hilfe des Zentralen Grenzwertsatzes erreicht werden.

Satz 5.1.1 (Zentraler Grenzwertsatz)

Seien $\mathbf{X}_k = (X_{k1}, \dots, X_{kd})'$, $k = 1, \dots, n$, $d < \infty$, unabhängig identisch verteilte Zufallsvektoren mit Erwartungswertvektor $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)'$ und existierender Kovarianzmatrix $\mathbf{V} < \infty$. Weiterhin sei \mathbf{T} ein Projektor (mit $H_0: \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ sowie $\mathbf{T}\mathbf{V}\mathbf{T} = \boldsymbol{\Sigma}$) und $\bar{\mathbf{X}}. = (\bar{X}_{.1}, \dots, \bar{X}_{.d})'$ der Vektor der arithmetischen Mittel. Dann gilt:

$$\begin{aligned} \sqrt{n} \cdot (\bar{\mathbf{X}}. - \boldsymbol{\mu}) &\xrightarrow{L} \mathbf{U} \sim N(\mathbf{0}, \mathbf{V}) \text{ und} \\ \sqrt{n} \cdot (\mathbf{T}\bar{\mathbf{X}}. - \mathbf{T}\boldsymbol{\mu}) &\xrightarrow{L} \mathbf{U} \sim N(\mathbf{0}, \boldsymbol{\Sigma}) \text{ für } n \rightarrow \infty. \end{aligned}$$

Beweis: Siehe zum Beispiel Ferguson (1996), S. 26-35.

□

Da die Zufallsvektoren $\mathbf{X}_k = \boldsymbol{\Gamma}\mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}$, $k = 1, \dots, n$, unter den Modellannahmen (2.1) die Voraussetzungen des Zentralen Grenzwertsatzes erfüllen, kann auf die quadratische Form der Mittelwertvektoren $Q_n^* = n \cdot \bar{\mathbf{X}}.'\mathbf{T}\bar{\mathbf{X}}.$ die Box-Approximation (4.2) asymptotisch angewendet werden und es folgt daraus:

$$\frac{Q_n^*}{Sp(\boldsymbol{\Sigma})} = \frac{n \cdot \bar{\mathbf{X}}.'\mathbf{T}\bar{\mathbf{X}}.}{Sp(\boldsymbol{\Sigma})} \dot{\sim} \chi_f^2/f, \quad f = \frac{[Sp(\boldsymbol{\Sigma})]^2}{Sp(\boldsymbol{\Sigma}^2)}. \quad (5.1)$$

5.2. Ansatz nach Geisser-Greenhouse

Zu Beginn wird analog zu Geisser-Greenhouse (4.4) die unbekannte Kovarianzmatrix im Nenner des Quotienten $Q_n^*/Sp(\boldsymbol{\Sigma})$ durch die empirische Kovarianzmatrix

$$\hat{\boldsymbol{\Sigma}}_n = \frac{1}{(n-1)} \sum_{k=1}^n (\mathbf{T}\mathbf{X}_k - \mathbf{T}\bar{\mathbf{X}}.)' (\mathbf{T}\mathbf{X}_k - \mathbf{T}\bar{\mathbf{X}}.)$$

ersetzt. Daraus ergibt sich der neue Quotient

$$C_n = \frac{Q_n^*}{Sp(\hat{\boldsymbol{\Sigma}}_n)}, \quad (5.2)$$

dessen Verteilung approximativ bestimmt werden muss.

Da unter Normalverteilung die beiden Zufallsvariablen im Zähler und Nenner unkorreliert sind, könnte dann eine weitere Box-Approximation auf den Nenner angewendet werden. Das ist allerdings unter den gegebenen Modellannahmen (2.1) nicht möglich. Denn aufgrund der Tatsache, dass unkorrelierte normalverteilte Linearformen stochastisch unabhängig sind, folgt auch die Unabhängigkeit der quadratischen Formen (siehe Lemma A.3.8 und Satz A.3.9 Craig und Sakamoto). Ohne die explizite Annahme der Normalverteilung folgt aber aus der Unkorreliertheit der Linearformen nicht mehr deren Unabhängigkeit, somit sind auch die quadratischen Formen nicht mehr unkorreliert. Es ist also nicht sinnvoll eine separate Box-Approximation für die Spur der empirischen Kovarianzmatrix im Nenner durchzuführen, da der Quotient zweier abhängiger χ^2 -Verteilungen nicht mehr F -verteilt wäre.

Das Problem der Abhängigkeit von Zähler und Nenner wird umgangen, indem der Quotient der beiden quadratischen Formen $C_n = Q_n^*/Sp(\widehat{\Sigma}_n)$ als eine neue Zufallsvariable aufgefasst wird.

Die Idee ist, die Verteilung der neuen Zufallsvariablen C_n mit Hilfe einer Box-Approximation anzunähern.

In (5.1) wurde bereits gezeigt, dass $\frac{Q_n^*}{Sp(\Sigma)}$ approximativ einer χ_f^2/f Verteilung folgt. Wird im Zähler die Spur der Kovarianzmatrix durch den erwartungstreuen und konsistenten Schätzer $\widehat{\Sigma}_n$ ersetzt, so ist die Verteilung des Quotienten $C_n = \frac{Q_n^*}{Sp(\widehat{\Sigma}_n)}$ unbekannt.

Da aber $\widehat{\Sigma}_n \xrightarrow{p} \Sigma$ konvergiert, kann C_n als eine Approximation des Quotienten $\frac{Q_n^*}{Sp(\Sigma)}$ angesehen werden.

Daher wird im Folgenden die Verteilung des Quotienten C_n mit Hilfe einer modifizierten Box-Approximation angenähert, so dass die ersten beiden Momente von C_n mit denen einer mit g_1 gestreckten $\chi_{f_1}^2/f_1$ -Verteilung übereinstimmen. Es wird also nicht mehr wie bisher, die Verteilung von Zähler und Nenner einzeln, durch zwei getrennte Box-Approximationen angenähert, sondern die Verteilung des Quotienten C_n wird approximiert.

In den folgenden Abschnitten soll die benötigte Theorie für diese Approximation hergeleitet werden.

5.3. Box-Approximation eines Quotienten von Zufallsvariablen

Es wird versucht über eine modifizierte Box-Approximation die Verteilung des Quotienten der quadratischen Formen C_n direkt durch eine mit g_1 gestreckte $\chi_{f_1}^2/f_1$ -Verteilung anzunähern, so dass Erwartungswert und Varianz übereinstimmen.

Sei dafür

$$U_1 \sim \chi_{f_1}^2/f_1$$

mit

$$E(C_n) = E\left(\frac{Q_n^*}{Sp(\widehat{\Sigma}_n)}\right) \stackrel{!}{=} E(g_1 U_1) = g_1,$$

$$Var(C_n) = Var\left(\frac{Q_n^*}{Sp(\widehat{\Sigma}_n)}\right) \stackrel{!}{=} Var(g_1 U_1) = 2 \cdot \frac{g_1^2}{f_1}.$$

Um die Approximation weiter fortsetzen zu können, werden der Erwartungswert und die Varianz des Quotienten C_n benötigt, welche nicht ohne weiteres exakt bestimmt werden können. Daher werden die beiden Momente durch die im Nachfolgenden erläuterte Taylor-Approximation angenähert.

Lemma 5.3.1 (Taylor-Approximation)

Seien X und Y zwei Zufallsvariablen mit $E(X) = \mu_x \neq 0$ und $E(Y) = \mu_y \neq 0$, sowie existierenden Varianzen $Var(X) = \sigma_x^2 < \infty$ und $Var(Y) = \sigma_y^2 < \infty$. Dann liefert die Taylor-Approximation für die ersten beiden Momente eines Quotienten zweier Zufallsvariablen folgendes:

$$E\left(\frac{X}{Y}\right) \doteq \frac{E(X)}{E(Y)} \cdot \left[1 + \frac{Var(Y)}{[E(Y)]^2} - \frac{Cov(X,Y)}{E(X) \cdot E(Y)}\right],$$

$$Var\left(\frac{X}{Y}\right) \doteq \frac{[E(X)]^2}{[E(Y)]^2} \cdot \left[\frac{Var(X)}{[E(X)]^2} + \frac{Var(Y)}{[E(Y)]^2} - 2 \cdot \frac{Cov(X,Y)}{E(X) \cdot E(Y)}\right].$$

Beweis: Siehe Strange (1970), S. 173-175.

□

Um die Momente der Zufallsvariable des Quotienten C_n mittels Taylor-Approximation annähern zu können, werden vorab die Erwartungswerte, Varianzen und die Kovarianz der quadratischen Formen im Zähler und Nenner benötigt. Für die Berechnung der Varianzen und Kovarianzen ist allerdings eine multivariate Version des Satzes von Atiqullah (1962) über die Varianz einer quadratischen Form erforderlich.

Dieser Satz existiert noch nicht in konkreter Form und wird im nächsten Abschnitt hergeleitet.

Es muss erwähnt werden, dass sich bereits Ohtaki (1990) mit der Darstellung der Kovarianz von quadratischen Formen beschäftigt hat. Dabei wurde allerdings ein Regressionsmodell der Form

$$\mathbf{Y} = \eta(\mathbf{X}) + \epsilon$$

angenommen, in dem die Parametervektoren \mathbf{Y} und \mathbf{X} gegeben waren, ϵ den Vektor der Fehlerterme bezeichnete und $\eta(\cdot)$ eine unbekannte Funktion war. Weiterhin wurde die dort beschriebene Herleitung der Kovarianz zweier quadratischer Formen mit Hilfe dieses Regressionsmodells durchgeführt. Da in der vorliegenden Arbeit keine Regressionsmodelle behandelt werden, müssen die benötigten Sätze über Varianzen und Kovarianzen von quadratischen Formen selbst hergeleitet werden.

Um die Herleitungen dieser Sätze zu vereinfachen, sei daran erinnert, dass sich für Zufallsvektoren $\mathbf{Y}_k = (Y_{k1}, \dots, Y_{kd})'$, $k = 1, \dots, n$, mit $\mathbf{Y} = (\mathbf{Y}'_1, \dots, \mathbf{Y}'_n)'$ und Matrix $\mathbf{A} = \mathbf{A}' = \mathbf{A}_n \otimes \mathbf{I}_d$ mit \mathbf{A}_n symmetrisch und Einträgen a_{ij} , $i, j = 1, \dots, n$, die quadratische Form $\mathbf{Y}'\mathbf{A}\mathbf{Y}$ wie folgt darstellen lässt:

$$\mathbf{Y}'\mathbf{A}\mathbf{Y} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} \cdot \mathbf{Y}'_i \mathbf{Y}_j.$$

5.3.1. Kovarianz zweier quadratischer Formen

Satz 5.3.2 (Kovarianz von quadratischen Formen)

Die Zufallsvektoren $\mathbf{Y}_k = (Y_{k1}, \dots, Y_{kd})'$, $k = 1, \dots, n$, seien unabhängig und identisch verteilt und bezeichne $\mathbf{Y} = (\mathbf{Y}'_1, \dots, \mathbf{Y}'_n)'$ den Vektor aller Zufallsvariablen mit $E_{H_0}(\mathbf{Y}) = \mathbf{0}$ und $Cov(\mathbf{Y}) = \mathbf{I}_n \otimes \Sigma$. Ferner sei $\tau_4 = E([\mathbf{Y}'_k \mathbf{Y}_k]^2)$, $k = 1, \dots, n$ und für alle $d < \infty$ soll $\tau_4 \leq \gamma_{\tau_4} < \infty$ gelten. Seien $\mathbf{A} = \mathbf{A}_n \otimes \mathbf{I}_d$ und $\mathbf{B} = \mathbf{B}_n \otimes \mathbf{I}_d$ symmetrische Matrizen mit $\mathbf{a}_n = \text{diag}\{\mathbf{A}_n\}$ und $\mathbf{b}_n = \text{diag}\{\mathbf{B}_n\}$.

Dann gilt:

$$\text{Cov}(\mathbf{Y}'\mathbf{A}\mathbf{Y}, \mathbf{Y}'\mathbf{B}\mathbf{Y}) = \left(\tau_4 - [\text{Sp}(\boldsymbol{\Sigma})]^2 - 2\text{Sp}(\boldsymbol{\Sigma}^2) \right) \mathbf{a}'_n \mathbf{b}_n + 2\text{Sp}(\boldsymbol{\Sigma}^2) \text{Sp}(\mathbf{A}_n \mathbf{B}_n).$$

Beweis: Siehe Anhang (Satz A.2.1) Seite 76

□

5.3.2. Varianz einer quadratischen Form

In Anlehnung an die vorangegangenen Überlegungen für den Satz über die Kovarianz zweier quadratischer Formen, lässt sich daraus leicht die Varianz einer quadratischen Form ableiten.

Satz 5.3.3 (Varianz einer quadratischen Form)

Die Zufallsvektoren $\mathbf{Y}_k = (Y_{k1}, \dots, Y_{kd})'$, $k = 1, \dots, n$, seien unabhängig und identisch verteilt und bezeichne $\mathbf{Y} = (\mathbf{Y}'_1, \dots, \mathbf{Y}'_n)'$ den Vektor aller Zufallsvariablen mit $E_{H_0}(\mathbf{Y}) = \mathbf{0}$ und $\text{Cov}(\mathbf{Y}) = \mathbf{I}_n \otimes \boldsymbol{\Sigma}$. Ferner sei $\tau_4 = E\left([\mathbf{Y}'_k \mathbf{Y}_k]^2\right)$, $k = 1, \dots, n$ und für alle $d < \infty$ soll $\tau_4 \leq \gamma_{\tau_4} < \infty$ gelten. Sei $\mathbf{A} = \mathbf{A}' = \mathbf{A}_n \otimes \mathbf{I}_d$ und bezeichne $\mathbf{a}_n = \text{diag}\{\mathbf{A}_n\}$.

Dann gilt:

$$\text{Var}(\mathbf{Y}'\mathbf{A}\mathbf{Y}) = \left(\tau_4 - [\text{Sp}(\boldsymbol{\Sigma})]^2 - 2\text{Sp}(\boldsymbol{\Sigma}^2) \right) \mathbf{a}'_n \mathbf{a}_n + 2\text{Sp}(\boldsymbol{\Sigma}^2) \text{Sp}(\mathbf{A}_n^2).$$

Beweis:

Folgt als Spezialfall aus dem Beweis von Satz 5.3.2 über die Kovarianz zweier quadratischer Formen, indem zwei identische quadratische Formen verwendet werden und somit $\mathbf{A} = \mathbf{A}_n \otimes \mathbf{I}_d = \mathbf{B}_n \otimes \mathbf{I}_d = \mathbf{B}$ gilt.

□

Für $d = 1$ ergibt sich die Formel für die Varianz einer quadratischen Form nach Atiqullah (1962) (siehe Satz A.3.7 Seite 89), somit kann dieser Satz als multivariate Version des Satzes von Atiqullah bezeichnet werden.

Im folgenden Abschnitt wird gezeigt, dass die Modellannahmen in (2.1) auf Seite 5 die Voraussetzungen der Sätze zur Darstellung der Varianz (Satz 5.3.3) und Kovarianz (Satz 5.3.2) von quadratischen Formen erfüllen.

5.4. Anwendung der Theorie

Da sich unter den Modellannahmen in (2.1) die Zufallsvektoren \mathbf{X}_k als

$$\mathbf{X}_k = \mathbf{\Gamma} \mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}, \quad k = 1, \dots, n,$$

darstellen lassen folgt aus Proposition 2.3.1 für $\mathbf{Y}_k = \mathbf{T} \mathbf{X}_k = \mathbf{T} \mathbf{\Gamma} \mathbf{Z}_k + \mathbf{T} \boldsymbol{\mu}$ unter $H_0: \mathbf{T} \boldsymbol{\mu} = \mathbf{0}$:

$$\begin{aligned} Y_k &= \mathbf{T} \mathbf{X}_k = \mathbf{T} \mathbf{\Gamma} \mathbf{Z}_k, \\ E_{H_0}(\mathbf{Y}_k) &= \mathbf{0}, \\ Cov(\mathbf{Y}_k) &= \mathbf{T} \mathbf{S} \mathbf{T} = \boldsymbol{\Sigma}, \quad k = 1, \dots, n. \end{aligned}$$

Sei dann $\mathbf{Y} = (\mathbf{Y}'_1, \dots, \mathbf{Y}'_n)'$ der Stichprobenvektor, mit $E(\mathbf{Y}) = \mathbf{0}$ und $Cov(\mathbf{Y}) = \mathbf{I}_n \otimes \boldsymbol{\Sigma}$. Es bleibt noch zu zeigen, dass $\tau_4 \leq \gamma_{\tau_4} < \infty \forall d < \infty$ gilt. Dafür werden Resultate der Lemmata 6.1.3 und 6.1.8 (Momente I und II) verwendet, welche in Kapitel 6 (Schätzer) hergeleitet werden, sowie die Modellannahme (2.1), dass die vierten Momente der Zufallsvariablen Z_{ks} beschränkt sind: $E(Z_{ks}^4) \leq \gamma < \infty \forall k = 1, \dots, n, s = 1, \dots, d$. Dann gilt für τ_4 :

$$\begin{aligned} \tau_4 &= E\left([\mathbf{Y}'_k \mathbf{Y}_k]^2\right) \stackrel{6.1.8}{\leq} \gamma Sp(\boldsymbol{\Sigma}^2) + [Sp(\boldsymbol{\Sigma})]^2 \stackrel{6.1.3}{\leq} (\gamma + 1) [Sp(\boldsymbol{\Sigma})]^2 \\ &\leq \gamma_{\tau_4} < \infty. \end{aligned}$$

Die letzte Ungleichung folgt aus der Tatsache, dass, solange die Dimension $d < \infty$ ist, die Spur der Kovarianzmatrix unter den Modellannahmen existiert und beschränkt ist mit:

$$[Sp(\boldsymbol{\Sigma})]^2 = \left[\sum_{i=1}^d \lambda_i \right]^2 \leq d^2 \cdot \left[\max_i (\lambda_i) \right] = d^2 \cdot \lambda_m^2 < \infty,$$

wobei die λ_i die Eigenwerte der Kovarianzmatrix $\boldsymbol{\Sigma}$ sind.

Das heißt, dass die Voraussetzungen für die Sätze 5.3.2 und 5.3.3 zur Berechnung der Varianzen und Kovarianzen von quadratischen Formen erfüllt sind und es möglich ist

die Sätze auf die quadratischen Formen des Quotienten C_n anzuwenden. Dazu werden die beiden Formen in der benötigten Matrizenschreibweise dargestellt.

Daraus folgt für die quadratische Form der Mittelwerte Q_n^* mit dem Vektor der arithmetischen Mittel $\bar{\mathbf{Y}}. = \left(\frac{1}{n}\mathbf{1}'_n \otimes \mathbf{I}_d\right) \mathbf{Y}$:

$$\begin{aligned} Q_n^* &= n \cdot \bar{\mathbf{X}}.' \mathbf{T} \bar{\mathbf{X}}. = n \cdot \bar{\mathbf{Y}}.' \bar{\mathbf{Y}}. = \mathbf{Y}' \left(\frac{1}{n} \mathbf{J}_n \otimes \mathbf{I}_d \right) \mathbf{Y} \\ &= \mathbf{Y}' (\mathbf{A}_n \otimes \mathbf{I}_d) \mathbf{Y} = \mathbf{Y}' \mathbf{A} \mathbf{Y} \end{aligned}$$

mit $\mathbf{A} = \mathbf{A}_n \otimes \mathbf{I}_d = \frac{1}{n} \mathbf{J}_n \otimes \mathbf{I}_d$.

Um die Spur der empirischen Stichprobenkovarianzmatrix als quadratische Form darstellen zu können, bedarf es einiger Vorbereitung.

Sei

$$(\mathbf{P}_n \otimes \mathbf{I}_d) \cdot \mathbf{Y} = \begin{pmatrix} \mathbf{Y}_1 - \bar{\mathbf{Y}}. \\ \vdots \\ \mathbf{Y}_n - \bar{\mathbf{Y}}. \end{pmatrix}$$

der mit $\bar{\mathbf{Y}}.$ zentrierte Zufallsvektor, dann stellt sich die Spur der empirischen Kovarianzmatrix wie folgt dar:

$$\begin{aligned} Sp(\hat{\Sigma}_n) &= Sp\left(\frac{1}{(n-1)} \sum_{k=1}^n (\mathbf{Y}_k - \bar{\mathbf{Y}}.) (\mathbf{Y}_k - \bar{\mathbf{Y}}.)'\right) \\ &= \frac{1}{(n-1)} \sum_{k=1}^n Sp\left((\mathbf{Y}_k - \bar{\mathbf{Y}}.) (\mathbf{Y}_k - \bar{\mathbf{Y}}.)'\right) \\ &= \frac{1}{(n-1)} \sum_{k=1}^n (\mathbf{Y}_k - \bar{\mathbf{Y}}.)' (\mathbf{Y}_k - \bar{\mathbf{Y}}.) \end{aligned}$$

Daraus ergibt sich die Möglichkeit $(n-1) \cdot Sp(\hat{\Sigma}_n)$ als quadratische Form darzustellen (siehe auch Geisser-Greenhouse (1968)):

$$\begin{aligned} (n-1) \cdot Sp(\hat{\Sigma}_n) &= \sum_{k=1}^n (\mathbf{Y}_k - \bar{\mathbf{Y}}.)' (\mathbf{Y}_k - \bar{\mathbf{Y}}.) = \mathbf{Y}' (\mathbf{P}_n \otimes \mathbf{I}_d) \mathbf{Y} \\ &= \mathbf{Y}' (\mathbf{B}_n \otimes \mathbf{I}_d) \mathbf{Y} = \mathbf{Y}' \mathbf{B} \mathbf{Y} \end{aligned}$$

mit $\mathbf{B} = \mathbf{B}_n \otimes \mathbf{I}_d = \mathbf{P}_n \otimes \mathbf{I}_d$.

Wird $\mathbf{B}^* = \mathbf{B}_n^* \otimes \mathbf{I}_d = \frac{1}{(n-1)} \cdot (\mathbf{P}_n \otimes \mathbf{I}_d)$ gewählt, so lässt sich $Sp(\hat{\Sigma}_n)$ direkt als quadratische Form darstellen: $Sp(\hat{\Sigma}_n) = \mathbf{Y}' \mathbf{B}^* \mathbf{Y}$.

Nachdem die quadratischen Formen in die gewünschte Gestalt gebracht wurden, können die Varianzen und die Kovarianz direkt berechnet werden, wobei im Nachfolgenden nur die Resultate aufgezeigt werden.

Die Varianz von Q_n^* ergibt sich wie folgt:

$$\begin{aligned} \text{Var}(Q_n^*) &= \text{Var}(\mathbf{Y}'\mathbf{A}\mathbf{Y}) = \text{Var}\left(\mathbf{Y}'\left(\frac{1}{n}\mathbf{J}_n \otimes \mathbf{I}_d\right)\mathbf{Y}\right) \\ &= \left(\tau_4 - [\text{Sp}(\boldsymbol{\Sigma})]^2 - 2\text{Sp}(\boldsymbol{\Sigma}^2)\right) \cdot \frac{1}{n} + 2\text{Sp}(\boldsymbol{\Sigma}^2). \end{aligned}$$

Analog lässt sich auch die Varianz von $\text{Sp}(\widehat{\boldsymbol{\Sigma}}_n)$ bestimmen:

$$\begin{aligned} \text{Var}\left(\text{Sp}\left(\widehat{\boldsymbol{\Sigma}}_n\right)\right) &= \text{Var}(\mathbf{Y}'\mathbf{B}^*\mathbf{Y}) = \frac{1}{(n-1)^2} \cdot \text{Var}(\mathbf{Y}'(\mathbf{P}_n \otimes \mathbf{I}_d)\mathbf{Y}) \\ &= \frac{1}{(n-1)^2} \left[\left(\tau_4 - [\text{Sp}(\boldsymbol{\Sigma})]^2 - 2\text{Sp}(\boldsymbol{\Sigma}^2)\right) \frac{(n-1)^2}{n} + 2\text{Sp}(\boldsymbol{\Sigma}^2)(n-1) \right] \\ &= \left(\tau_4 - [\text{Sp}(\boldsymbol{\Sigma})]^2 - 2\text{Sp}(\boldsymbol{\Sigma}^2)\right) \frac{1}{n} + 2\text{Sp}(\boldsymbol{\Sigma}^2) \frac{1}{(n-1)}. \end{aligned}$$

Abschließend wird die Kovarianz der quadratischen Formen Q_n^* und $\text{Sp}(\widehat{\boldsymbol{\Sigma}}_n)$ berechnet:

$$\begin{aligned} \text{Cov}\left(Q_n^*, \text{Sp}\left(\widehat{\boldsymbol{\Sigma}}_n\right)\right) &= \text{Cov}(\mathbf{Y}'\mathbf{A}\mathbf{Y}, \mathbf{Y}'\mathbf{B}^*\mathbf{Y}) = \frac{1}{(n-1)} \text{Cov}(\mathbf{Y}'\mathbf{A}\mathbf{Y}, \mathbf{Y}'\mathbf{B}\mathbf{Y}) \\ &= \frac{1}{(n-1)} \left[\left(\tau_4 - [\text{Sp}(\boldsymbol{\Sigma})]^2 - 2\text{Sp}(\boldsymbol{\Sigma}^2)\right) \frac{(n-1)}{n} + 2\text{Sp}(\boldsymbol{\Sigma}^2) \cdot 0 \right] \\ &= \frac{1}{n} \cdot \left(\tau_4 - [\text{Sp}(\boldsymbol{\Sigma})]^2 - 2 \cdot \text{Sp}(\boldsymbol{\Sigma}^2)\right). \end{aligned}$$

5.4.1. Taylor-Approximation

Sind die Varianzen und Kovarianzen der quadratischen Formen bestimmt, lässt sich die Approximation der Momente des Quotienten C_n mittels der in Lemma 5.3.1 definierten Taylor-Approximation durchführen.

Dafür werden die Ergebnisse der vorangegangenen Kapitel kurz aufgelistet.

Sei weiterhin

$$C_n = \frac{n \cdot \bar{\mathbf{Y}}' \bar{\mathbf{Y}}}{Sp(\hat{\Sigma}_n)} = \frac{Q_n^*}{Sp(\hat{\Sigma}_n)}$$

mit den zugehörigen Erwartungswerten, Varianzen und Kovarianzen der quadratischen Formen Q_n^* und $\hat{\Sigma}_n$:

$$\begin{aligned} E(Q_n^*) &= Sp(\Sigma), \\ E(Sp(\hat{\Sigma}_n)) &= Sp(\Sigma), \\ Var(Q_n^*) &= \frac{1}{n} (\tau_4 - [Sp(\Sigma)]^2 - 2Sp(\Sigma^2)) + 2Sp(\Sigma^2), \\ Var(Sp(\hat{\Sigma}_n)) &= \frac{1}{n} (\tau_4 - [Sp(\Sigma)]^2 - 2Sp(\Sigma^2)) + 2Sp(\Sigma^2) \frac{1}{(n-1)}, \\ Cov(Q_n^*, Sp(\hat{\Sigma}_n)) &= \frac{1}{n} (\tau_4 - [Sp(\Sigma)]^2 - 2Sp(\Sigma^2)). \end{aligned}$$

Zur Vereinfachung der Schreibweise werden folgende Bezeichnungen verwendet:

$$\beta_1 = [Sp(\Sigma)]^2 \text{ und } \beta_2 = Sp(\Sigma^2).$$

Somit sind die Voraussetzungen für Lemma 5.3.1 erfüllt und der Erwartungswert und die Varianz des Quotienten $C_n = \frac{Q_n^*}{Sp(\hat{\Sigma}_n)}$ können angenähert werden:

$$\begin{aligned} E\left(\frac{Q_n^*}{Sp(\hat{\Sigma}_n)}\right) &\doteq \frac{E(Q_n^*)}{E(Sp(\hat{\Sigma}_n))} \left[1 + \frac{Var(Sp(\hat{\Sigma}_n))}{[E(Sp(\hat{\Sigma}_n))]^2} - \frac{Cov(Q_n^*, Sp(\hat{\Sigma}_n))}{E(Q_n^*) E(Sp(\hat{\Sigma}_n))} \right] \\ &= \frac{Sp(\Sigma)}{Sp(\Sigma)} \left[1 + \frac{\frac{1}{n} (\tau_4 - \beta_1 - 2\beta_2) + 2\beta_2 \frac{1}{(n-1)}}{\beta_1} - \frac{\frac{1}{n} (\tau_4 - \beta_1 - 2\beta_2)}{\beta_1} \right] \\ &= \frac{\beta_1 + 2\beta_2 \cdot \frac{1}{(n-1)}}{\beta_1} = \frac{[Sp(\Sigma)]^2 + 2Sp(\Sigma^2) \cdot \frac{1}{(n-1)}}{[Sp(\Sigma)]^2}, \end{aligned}$$

$$\begin{aligned}
 \text{Var} \left(\frac{Q_n^*}{Sp(\widehat{\Sigma}_n)} \right) &\doteq \frac{[E(Q_n^*)]^2}{[E(Sp(\widehat{\Sigma}_n))]^2} \left[\frac{\text{Var}(Q_n^*)}{[E(Q_n^*)]^2} + \frac{\text{Var}(Sp(\widehat{\Sigma}_n))}{[E(Sp(\widehat{\Sigma}_n))]^2} \right] \\
 &\quad - \frac{[E(Q_n^*)]^2}{[E(Sp(\widehat{\Sigma}_n))]^2} \left[2 \frac{\text{Cov}(Q_n^*, Sp(\widehat{\Sigma}_n))}{E(Q_n^*) E(Sp(\widehat{\Sigma}_n))} \right] \\
 &= \frac{2Sp(\Sigma^2) \cdot \frac{n}{(n-1)}}{[Sp(\Sigma)]^2}.
 \end{aligned}$$

5.4.2. Box-Approximation

Die Ergebnisse der Taylor-Approximation werden benutzt, um die in Abschnitt 5.3 begonnene Box-Approximation des Quotienten $C_n = \frac{Q_n^*}{Sp(\widehat{\Sigma}_n)}$ berechnen zu können.

Sei dafür:

$$U_1 \sim \chi_{f_1}^2 / f_1$$

mit

$$E(C_n) = E \left(\frac{Q_n^*}{Sp(\widehat{\Sigma}_n)} \right) \doteq \frac{[Sp(\Sigma)]^2 + 2 \cdot Sp(\Sigma^2) \cdot \frac{1}{(n-1)}}{[Sp(\Sigma)]^2} = E(g_1 U_1) = g_1,$$

$$\text{Var}(C_n) = \text{Var} \left(\frac{Q_n^*}{Sp(\widehat{\Sigma}_n)} \right) \doteq \frac{2 \cdot Sp(\Sigma^2) \cdot \frac{n}{(n-1)}}{[Sp(\Sigma)]^2} = \text{Var}(g_1 U_1) = 2 \cdot \frac{g_1^2}{f_1}.$$

Der Streckungsparameter g_1 kann aus diesem Gleichungssystem direkt abgelesen werden:

$$\begin{aligned}
 g_1 &= \frac{[Sp(\Sigma)]^2 + 2 \cdot Sp(\Sigma^2) \cdot \frac{1}{(n-1)}}{[Sp(\Sigma)]^2} \\
 &= 1 + \frac{2}{(n-1)} \cdot \frac{Sp(\Sigma^2)}{[Sp(\Sigma)]^2} = 1 + \frac{2}{(n-1) \cdot f}
 \end{aligned}$$

mit $f = \frac{[Sp(\boldsymbol{\Sigma})]^2}{Sp(\boldsymbol{\Sigma}^2)}$, dem Freiheitsgrad aus der Box-Approximation in (5.1) (mit bekannter Kovarianzmatrix $\boldsymbol{\Sigma}$). Für $n \rightarrow \infty$ geht der Streckungsparameter $g_1 \rightarrow 1$.

Danach wird g_1 in die zweite Gleichung eingesetzt und diese nach f_1 aufgelöst:

$$\begin{aligned} f_1 &= \frac{2 \cdot g_1^2}{Var(g_1 \cdot U_1)} \\ &= 2 \cdot \left[\frac{[Sp(\boldsymbol{\Sigma})]^2 + 2 \cdot Sp(\boldsymbol{\Sigma}^2) \cdot \frac{1}{(n-1)}}{[Sp(\boldsymbol{\Sigma})]^2} \right]^2 \cdot \frac{[Sp(\boldsymbol{\Sigma})]^2}{2 \cdot Sp(\boldsymbol{\Sigma}^2) \cdot \frac{n}{(n-1)}} \\ &= \frac{[(n-1) \cdot [Sp(\boldsymbol{\Sigma})]^2 + 2 \cdot Sp(\boldsymbol{\Sigma}^2)]^2}{[Sp(\boldsymbol{\Sigma})]^2 \cdot Sp(\boldsymbol{\Sigma}^2) \cdot n(n-1)} \\ &= f \cdot \frac{(n-1)}{n} + \frac{4}{n} + \frac{4}{f \cdot n(n-1)}, \end{aligned}$$

wobei $f = \frac{[Sp(\boldsymbol{\Sigma})]^2}{Sp(\boldsymbol{\Sigma}^2)}$ analog zur Bestimmung von g_1 definiert ist und für $n \rightarrow \infty$ der Freiheitsgrad $f_1 \rightarrow f$ geht.

Nun ergibt sich die Möglichkeit, die Statistik C_n durch eine mit g_1 gestreckte $\chi_{f_1}^2/f_1$ -Verteilung zu approximieren, wobei für $n \rightarrow \infty$ C_n approximativ einer χ_f^2/f -Verteilung (5.1) folgt.

$$\frac{C_n}{g_1} = \frac{Q_n^*}{Sp(\hat{\boldsymbol{\Sigma}}_n) \cdot g_1} \dot{\sim} \chi_{f_1}^2/f_1 \quad (5.3)$$

Allerdings enthalten sowohl f_1 als auch g_1 die unbekannt Parameter $[Sp(\boldsymbol{\Sigma})]^2$ und $Sp(\boldsymbol{\Sigma}^2)$, die im Folgenden erwartungstreu, konsistent und dimensionsstabil geschätzt werden müssen. Da die Herleitung technisch sehr aufwendig ist, wird der Nachweis der gewünschten Eigenschaften für die Schätzer von $[Sp(\boldsymbol{\Sigma})]^2$ und $Sp(\boldsymbol{\Sigma}^2)$ in Kapitel 6 behandelt.

5.5. Die neue ANOVA-Typ Statistik

Die in Kapitel 6 hergeleiteten Schätzer B_1 und B_2 für $[Sp(\Sigma)]^2$ und $Sp(\Sigma^2)$ erfüllen alle gewünschten Eigenschaften und es ist möglich, den Freiheitsgrad f_1 und den Streckungsparameter g_1 durch Einsetzen von B_1 und B_2 zu schätzen.

$$\hat{f}_1 = \frac{[(n-1) \cdot B_1 + 2 \cdot B_2]^2}{B_1 \cdot B_2 \cdot n(n-1)}$$

$$\hat{g}_1 = 1 + \frac{2 \cdot B_2}{(n-1) \cdot B_1}$$

Durch Einsetzen von \hat{g}_1 und \hat{f}_1 in (5.3) ergibt sich eine neue Teststatistik, welche im Folgenden mit ATS-neu bezeichnet werden soll:

$$A_n^{ATS-neu} = \frac{C_n}{\hat{g}_1} = \frac{Q_n^*}{Sp(\hat{\Sigma}_n) \cdot \hat{g}_1} \dot{\sim} \chi_{\hat{f}_1}^2 / \hat{f}_1. \quad (5.4)$$

Diese Statistik wurde ohne explizite Annahme der Normalverteilung hergeleitet und kann daher auf beliebige metrische Daten, welche die Modellannahmen (Abschnitt 2.2.1) erfüllen, angewendet werden. Einzelheiten des Verhaltens der Teststatistik unter verschiedenen Verteilungen finden sich im Kapitel 8 (Simulationen).

Anmerkung

Wie am Ergebnis ersichtlich wird, enthalten die Parameter f und g das τ_4 nicht mehr, weshalb an dieser Stelle auf eine Definition und Herleitung eines Schätzers für τ_4 verzichtet wird.

Im folgenden Kapitel werden die bereits aus Werner (2004) bekannten Schätzer der Spuren der Kovarianzmatrix auf Nicht-Normalverteilung erweitert.

6. Die Schätzer B_1 und B_2

6.1. Quadrat- und Bilinearformen

Allgemein seien die Quadrat- und Bilinearformen analog zu Werner (2004) sowie Ahmad et al. (2008) definiert:

Definition 6.1.1 *Es gilt folgende Unterscheidung:*

$$A_{kl} = \mathbf{X}'_k \mathbf{T} \mathbf{X}_l = \mathbf{Y}'_k \mathbf{Y}_l \hat{=} \begin{cases} \text{Quadratform} & k = l \\ \text{Bilinearform} & k \neq l, \end{cases}$$

wobei für $k = l$ A_{kl} mit A_k bezeichnet werden soll.

Definition 6.1.2 *Seien B_1 und B_2 Schätzer für $[Sp(\boldsymbol{\Sigma})]^2$ und $Sp(\boldsymbol{\Sigma}^2)$ und analog zu Ahmad et al. (2008) wie folgt definiert:*

$$B_1 = \frac{1}{n(n-1)} \cdot \sum_{k \neq l} A_k \cdot A_l \qquad B_2 = \frac{1}{n \cdot (n-1)} \cdot \underbrace{\sum_{k=1}^n \sum_{l=1}^n}_{k \neq l} A_{kl}^2$$

Es folgen einige Lemmata und Sätze, die für den Nachweis der gewünschten Eigenschaften dieser Schätzer benötigt werden.

Lemma 6.1.3 (Momente I)

Seien $X_k = (X_{k1}, \dots, X_{kd})'$, $k = 1, \dots, n$, unabhängig identisch verteilte Vektoren von Zufallsvariablen mit $E(\mathbf{X}_k) = \boldsymbol{\mu}$, $Cov(\mathbf{X}_k) = \mathbf{S}$ und sei \mathbf{T} ein Projektor für den $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ gilt, sowie $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k$ mit $E_{H_0}(\mathbf{Y}_k) = \mathbf{0}$ und $Cov(\mathbf{Y}_k) = \mathbf{T}\mathbf{S}\mathbf{T} = \boldsymbol{\Sigma}$. Dann gilt für die Momente der Quadrat- und Bilinearformen:

1. $E(A_k) = Sp(\boldsymbol{\Sigma})$
2. $E(A_k A_l) \stackrel{k \neq l}{=} [Sp(\boldsymbol{\Sigma})]^2$
3. $E(A_{kl}) = 0$
4. $E(A_{kl}^2) = Sp(\boldsymbol{\Sigma}^2)$
5. $Sp(\boldsymbol{\Sigma}^2) \leq [Sp(\boldsymbol{\Sigma})]^2 \leq E(A_k^2)$

Beweis:

1.) folgt aus dem Satz von Lancaster (A.3.6)

$$E(A_k) = E(\mathbf{Y}'_k \mathbf{Y}_k) = Sp(\boldsymbol{\Sigma})$$

2.) folgt aus 1.) und der Unabhängigkeit der quadratischen Formen, für $k \neq l$.

3.) folgt ebenfalls aus der Unabhängigkeit, für $k \neq l$:

$$E(A_{kl}) = E(\mathbf{Y}'_k \mathbf{Y}_l) = E(\mathbf{Y}'_k) \cdot E(\mathbf{Y}_l) = 0.$$

4.) folgt durch Ausnutzen der Invarianz der Spur unter zyklischen Vertauschungen:

$$\begin{aligned} E(A_{kl}^2) &= E(\mathbf{Y}'_k \mathbf{Y}_l \mathbf{Y}'_k \mathbf{Y}_l) = E(Sp(\mathbf{Y}'_k \mathbf{Y}_l \mathbf{Y}'_k \mathbf{Y}_l)) \\ &= E(Sp(\mathbf{Y}'_k \mathbf{Y}_k \mathbf{Y}'_l \mathbf{Y}_l)) = Sp(E(\mathbf{Y}'_k \mathbf{Y}_k) \cdot E(\mathbf{Y}'_l \mathbf{Y}_l)) = Sp(\boldsymbol{\Sigma}^2). \end{aligned}$$

5.) Der erste Teil der Ungleichung folgt aus der Cauchy-Schwarz Ungleichung und der Eigenschaft, dass die Spur einer symmetrischen Matrix gleich der Summe ihrer Eigenwerte ist, wobei alle Eigenwerte $\lambda_i \geq 0$ sind da $\boldsymbol{\Sigma}$ positiv semidefinit ist. Der zweite Teil folgt aus der Definition der Varianz und 2.)

$$\begin{aligned} Sp(\boldsymbol{\Sigma}^2) &= \sum_{i=1}^d \lambda_i^2 \leq \left(\sum_{i=1}^d \lambda_i \right)^2 = [Sp(\boldsymbol{\Sigma})]^2 \\ Var(A_k) &= E(A_k^2) - [E(A_k)]^2 \geq 0 \\ \Rightarrow [Sp(\boldsymbol{\Sigma})]^2 &\leq E(A_k^2). \end{aligned}$$

□

Für die Berechnung und Abschätzung der höheren Momente der Quadrat- und Bilinearformen werden zuvor zwei spezielle Darstellungssätze benötigt, die im Folgenden hergeleitet werden.

6.1.1. Darstellungssätze

Satz 6.1.4 (Darstellung einer Bilinearform)

Seien $\mathbf{X} = (X_1, \dots, X_d)'$ und $\mathbf{Y} = (Y_1, \dots, Y_d)'$ unabhängig identisch verteilte Zufallsvektoren mit $E(\mathbf{X}) = E(\mathbf{Y}) = \mathbf{0}$ und $Cov(\mathbf{X}) = Cov(\mathbf{Y}) = \mathbf{S}$ mit $r(\mathbf{S}) = r \leq d$ sowie \mathbf{T} ein Projektor. Dann gilt für die Bilinearform:

$$A = \mathbf{X}'\mathbf{T}\mathbf{Y} = \sum_{i=1}^d \lambda_i U_i W_i,$$

wobei die λ_i die Eigenwerte von $\mathbf{T}\mathbf{S}$ sind. Außerdem sind $\mathbf{U} = (U_1, \dots, U_d)'$ und $\mathbf{W} = (W_1, \dots, W_d)'$ Zufallsvektoren mit $E(\mathbf{U}) = E(\mathbf{W}) = \mathbf{0}$ und $Cov(\mathbf{U}) = Cov(\mathbf{W}) = \mathbf{I}_d$.

Beweis:

1. Schritt

Die Aussage wird zunächst für $r = d$ bewiesen, d.h. $\det(\mathbf{S}) = |\mathbf{S}| \neq 0$. Dieser Fall wurde bereits von Werner (2004) gezeigt.

2. Schritt

Sei \mathbf{S} singulär, mit $r(\mathbf{S}) = r < d$.

Beweis:

Folgt aus dem Anwenden der Ergebnisse aus dem 1. Schritt auf den singulären Fall des Beweises für quadratische Formen von Mathai und Provost (1992), Seite 28-38 sowie Brunner (2010) Seite 36-38.

□

Korollar 6.1.5 (Darstellung einer Bilinearform)

Für die Zufallsvektoren gelte $\mathbf{X}_k = (X_{k1}, \dots, X_{kd})' = \mathbf{\Gamma}\mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}$, $k = 1, \dots, n$, wie in (2.1) mit $\text{Cov}(\mathbf{\Gamma}\mathbf{Z}_k) = \mathbf{\Gamma}\mathbf{\Gamma}' = \mathbf{S}$ und sei $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k + \mathbf{T}\boldsymbol{\mu}$ definiert wie in Proposition 2.3.1 mit $\text{Cov}(\mathbf{Y}_k) = \mathbf{T}\mathbf{S}\mathbf{T} = \boldsymbol{\Sigma}$. Dann ergibt sich für $k \neq l$ unter $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ das spezielle Resultat von Satz 6.1.4:

$$\begin{aligned} A_{kl} &= \mathbf{X}'_k \mathbf{T} \mathbf{X}_l = \mathbf{Y}'_k \mathbf{Y}_l \\ &= \sum_{i=1}^d \lambda_i Z_{ki} Z_{li}, \end{aligned}$$

wobei die λ_i die Eigenwerte von $\mathbf{T}\mathbf{S}$ sind und die Zufallsvariablen Z_{ki} unabhängig $\forall k = 1, \dots, n, i = 1, \dots, d$.

Beweis: Siehe Anhang (Satz A.2.2) Seite 80

□

Satz 6.1.6 (Darstellung einer Quadratform)

Sei $\mathbf{X} = (X_1, \dots, X_d)'$ ein Zufallsvektor mit $E(\mathbf{X}) = \mathbf{0}$ und $\text{Cov}(\mathbf{X}) = \mathbf{S}$ mit $r(\mathbf{S}) = r \leq d$ sowie \mathbf{T} ein Projektor. Dann gilt für die quadratische Form:

$$A = \mathbf{X}' \mathbf{T} \mathbf{X} = \sum_{i=1}^d \lambda_i U_i^2,$$

wobei die λ_i die Eigenwerte von $\mathbf{T}\mathbf{S}$ sind und $\mathbf{U} = (U_1, \dots, U_d)'$ ein Zufallsvektor mit $E(\mathbf{U}) = \mathbf{0}$ und $\text{Cov}(\mathbf{U}) = \mathbf{I}_d$ ist.

Beweis: Siehe Mathai und Provost (1992), Seite 28-38 sowie Werner (2004)

□

Korollar 6.1.7 (Darstellung einer Quadratform)

Für die Zufallsvektoren gelte $\mathbf{X}_k = (X_{k1}, \dots, X_{kd})' = \mathbf{\Gamma} \mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}$, $k = 1, \dots, n$, wie in (2.1) mit $\text{Cov}(\mathbf{\Gamma} \mathbf{Z}_k) = \mathbf{\Gamma} \mathbf{\Gamma}' = \mathbf{S}$ und sei $\mathbf{Y}_k = \mathbf{T} \mathbf{X}_k = \mathbf{T} \mathbf{\Gamma} \mathbf{Z}_k + \mathbf{T} \boldsymbol{\mu}$ definiert wie in Proposition 2.3.1 mit $\text{Cov}(\mathbf{Y}_k) = \mathbf{T} \mathbf{S} \mathbf{T}' = \boldsymbol{\Sigma}$. Daraus ergibt sich unter $H_0 : \mathbf{T} \boldsymbol{\mu} = \mathbf{0}$ eine speziellere Version von Satz 6.1.6 und für die quadratische Form gilt:

$$\begin{aligned} A_k &= \mathbf{X}_k' \mathbf{T} \mathbf{X}_k = \mathbf{Y}_k' \mathbf{Y}_k \\ &= \sum_{i=1}^d \lambda_i Z_{ki}^2, \end{aligned}$$

wobei die λ_i die Eigenwerte von $\mathbf{T} \mathbf{S}$ sind und die Zufallsvariablen Z_{ki} unabhängig $\forall k = 1, \dots, n, i = 1, \dots, d$.

Beweis:

Ergibt sich als Spezialfall aus dem Beweis des Darstellungssatzes für Bilinearformen (Satz 6.1.5), indem der Zufallsvektor $\mathbf{X}_l = \mathbf{X}_k$ gesetzt wird.

□

Im Folgenden werden alle Berechnungen für \mathbf{S} positiv definit und mit existierender Wurzel durchgeführt, da die Rechnungen im singulären Fall analog verlaufen und sich nur der Index von d auf r reduziert (siehe Satz A.2.2).

Lemma 6.1.8 (Momente II)

Für die unabhängig identisch verteilten Zufallsvektoren gelte $\mathbf{X}_k = (X_{k1}, \dots, X_{kd})' = \mathbf{\Gamma}\mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}$, $k = 1, \dots, n$, wie in (2.1) mit $\text{Cov}(\mathbf{\Gamma}\mathbf{Z}_k) = \mathbf{\Gamma}\mathbf{\Gamma}' = \mathbf{S}$ und $E(Z_{ks}^4) \leq \gamma < \infty \forall k = 1, \dots, n, s = 1, \dots, d$. Sei $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k + \mathbf{T}\boldsymbol{\mu}$ definiert wie in Proposition 2.3.1 mit $\text{Cov}(\mathbf{Y}_k) = \mathbf{T}\mathbf{S}\mathbf{T}' = \boldsymbol{\Sigma}$. Des Weiteren seien $A_k = \mathbf{Y}_k' \mathbf{Y}_k$ und $A_{kl} = \mathbf{Y}_k' \mathbf{Y}_l$ $k \neq l$ Quadrat- und Bilinearformen (Siehe Definition 6.1.1).

Dann gilt unter $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ für die Momente der Formen:

- 1.) $E(A_k^2) \leq \gamma \text{Sp}(\boldsymbol{\Sigma}^2) + [\text{Sp}(\boldsymbol{\Sigma})]^2$
- 2.) $E(A_k^2 A_l^2) \leq \left[\gamma \text{Sp}(\boldsymbol{\Sigma}^2) + [\text{Sp}(\boldsymbol{\Sigma})]^2 \right]^2$
- 3.) $\text{Var}(A_k) \leq \gamma \text{Sp}(\boldsymbol{\Sigma}^2)$
- 4.) $E(A_{kl}^4) \leq \gamma \text{Sp}(\boldsymbol{\Sigma}^4) + 3 [\text{Sp}(\boldsymbol{\Sigma}^2)]^2$
- 5.) $E(A_{kl}^2 A_{rs}^2) = \begin{cases} E(A_{kl}^4) & \text{für } k = r, l = s \text{ und } k = s, l = r \\ [\text{Sp}(\boldsymbol{\Sigma}^2)]^2 & \text{für } k \neq l \neq r \neq s \end{cases}$

$$E(A_{kl}^2 A_{rs}^2) \leq E(A_{kl}^4) \quad \text{sonst.}$$

Beweis:

Aus Proposition 2.3.1 folgt unter $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$:

$\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k$ sowie $E_{H_0}(\mathbf{Y}_k) = \mathbf{0}$. Weiterhin sind die Zufallsvektoren \mathbf{Z}_k , $k = 1, \dots, n$, nach den Modellannahmen in (2.1) auf Seite 5 unabhängig und enthalten unabhängige Komponenten Z_{ki} , $i = 1, \dots, d$, mit beschränkten vierten Momenten: $E(Z_{ks}^4) \leq \gamma < \infty \forall k = 1, \dots, n, i = 1, \dots, d$.

1.) Da die quadratischen Formen die Voraussetzungen von Korollar 6.1.7 erfüllen, lässt sich $E(A_k^2)$ mit Hilfe der Zufallsvariablen Z_{ki} , $k = 1, \dots, n$, $i = 1, \dots, d$, wie folgt darstellen:

$$\begin{aligned} E(A_k^2) &= E \left[\left(\sum_{i=1}^d \lambda_i Z_{ki}^2 \right)^2 \right] = \sum_{i=1}^d \sum_{j=1}^d \lambda_i \lambda_j E(Z_{ki}^2 Z_{kj}^2) \\ &= \sum_{i=1}^d \lambda_i^2 E(Z_{ki}^4) + \sum_{i \neq j} \lambda_i \lambda_j E(Z_{ki}^2) E(Z_{kj}^2) \\ &\leq \gamma \cdot \sum_{i=1}^d \lambda_i^2 + 1 \cdot \sum_{i \neq j} \lambda_i \lambda_j \leq \gamma \text{Sp}(\boldsymbol{\Sigma}^2) + \sum_{i=1}^d \sum_{j=1}^d \lambda_i \lambda_j \end{aligned}$$

$$\begin{aligned}
 &= \gamma Sp(\Sigma^2) + \left(\sum_{i=1}^d \lambda_i \right) \left(\sum_{j=1}^d \lambda_j \right) \\
 &= \gamma Sp(\Sigma^2) + [Sp(\Sigma)]^2.
 \end{aligned}$$

2.) folgt aus 1.) und aus der Unabhängigkeit der quadratischen Formen für $k \neq l$:

$$\begin{aligned}
 E(A_k^2 A_l^2) &= E(A_k^2) E(A_l^2) = [E(A_k^2)]^2 \\
 &\leq \left[\gamma Sp(\Sigma^2) + [Sp(\Sigma)]^2 \right]^2.
 \end{aligned}$$

3.) folgt direkt aus 1.) und der Definition der Varianz.

4.) Da die Bilinearform A_{kl} die Voraussetzungen von Korollar 6.1.5 erfüllt, lässt sich $E(A_{kl}^4)$ ebenfalls über die Z_{ki} darstellen:

$$\begin{aligned}
 E(A_{kl}^4) &= E(A_{kl} A_{kl} A_{kl} A_{kl}) \\
 &= E \left[\sum_{i=1}^d \sum_{j=1}^d \sum_{r=1}^d \sum_{s=1}^d \lambda_i \lambda_j \lambda_r \lambda_s Z_{ki} Z_{li} Z_{kj} Z_{lj} Z_{kr} Z_{lr} Z_{ks} Z_{ls} \right] \\
 &= \sum_{i=1}^d \sum_{j=1}^d \sum_{r=1}^d \sum_{s=1}^d \lambda_i \lambda_j \lambda_r \lambda_s [E(Z_{ki} Z_{kj} Z_{kr} Z_{ks})]^2
 \end{aligned}$$

und aus der Unabhängigkeit der Z_{ki} , $k = 1, \dots, n$, $i = 1, \dots, d$, (siehe (2.1)) folgt:

$$\begin{aligned}
 E(A_{kl}^4) &= \sum_{i=1}^d \lambda_i^4 [E(Z_{ki}^4)]^2 + 3 \sum_{i \neq r} \lambda_i^2 \lambda_r^2 [E(Z_{ki}^2)]^4 \\
 &\leq \gamma \sum_{i=1}^d \lambda_i^4 + 3 \sum_{i \neq r} \lambda_i^2 \lambda_r^2 \leq \gamma \sum_{i=1}^d \lambda_i^4 + 3 \sum_{i=1}^d \sum_{r=1}^d \lambda_i^2 \lambda_r^2 \\
 &= \gamma Sp(\Sigma^4) + 3 \left(\sum_{i=1}^d \lambda_i^2 \right) \left(\sum_{r=1}^d \lambda_r^2 \right) = \gamma Sp(\Sigma^4) + 3 [Sp(\Sigma^2)]^2.
 \end{aligned}$$

5.) folgt aus der Unabhängigkeit, für $k \neq l$, und aus 4.):

I)

Für $k = r$, $l = s$ und analog für $k = s$, $r = l$ gilt

$$E(A_{kl}^2 \cdot A_{rs}^2) = E(A_{kl}^4)$$

da $E(A_{kl}) = E(A_{lk})$.

II)

Für alle Indizes verschieden folgt aus der Unabhängigkeit

$$E(A_{kl}^2 \cdot A_{rs}^2) = E(A_{kl}^2) \cdot E(A_{rs}^2) = [Sp(\Sigma^2)]^2.$$

III)

Für $k = r, l \neq s$ oder $k = s, l \neq r$ und analog

für $k \neq r, l = s$ oder $k \neq s, l = r$ folgt mit Hilfe der Cauchy-Schwarz Ungleichung:

$$E(A_{kl}^2 \cdot A_{ks}^2) \stackrel{CS}{\leq} \sqrt{E(A_{kl}^4) \cdot E(A_{ks}^4)} = E(A_{kl}^4).$$

□

Damit sind alle erforderlichen Vorbereitungen für den Nachweis der gewünschten Schätzeigenschaften getroffen und es kann der nachfolgende Satz hergeleitet werden.

6.2. Nachweis der Schätzeigenschaften

Satz 6.2.1 (Schätzeigenschaften)

Seien B_1 und B_2 definiert wie in Definition 6.1.2 und erfüllen die Zufallsvektoren \mathbf{X}_k , $k = 1, \dots, n$, die Modellannahmen aus (2.1). Weiterhin sei $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k + \mathbf{T}\boldsymbol{\mu}$ definiert wie in Proposition 2.3.1 mit $\text{Cov}(\mathbf{Y}_k) = \mathbf{T}\mathbf{S}\mathbf{T} = \boldsymbol{\Sigma}$.

Somit gilt unter $H_0: \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$: B_1 und B_2

1. sind erwartungstreue Schätzer für $[\text{Sp}(\boldsymbol{\Sigma})]^2$ und $\text{Sp}(\boldsymbol{\Sigma}^2)$.
2. sind konsistent im Sinne von Definition A.3.2.
3. sind dimensionsstabil im Sinne von Definition A.3.4.

Beweis: Siehe Anhang (A.2.3) Seite 84

□

In den vorangegangenen Kapiteln wurde eine neue Prüfgröße zum Testen der Hypothese $H_0: \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ ohne Annahme der Normalverteilung der Zufallsvariablen \mathbf{X}_k hergeleitet. Es wurden lediglich die Modellannahmen aus (2.1) vorausgesetzt. Für die gesamten Herleitungen wurde die Dimension der Messwiederholungen als beliebig, aber fest mit $d < \infty$ angenommen. Um allerdings die Konsistenz der Schätzer sowie die asymptotischen Resultate aus dem Zentralen Grenzwertsatz (Abschnitt 5.1.1) zu erhalten, ging der Stichprobenumfang $n \rightarrow \infty$.

Im nächsten Abschnitt soll die Statistik A_n^{ATS-W} (4.5) von Werner (2004) ohne die Annahme der Normalverteilung hergeleitet werden. Es wird sich zeigen, dass die gleiche Prüfgröße zum Testen von H_0 auch bei Daten eingesetzt werden kann, die keiner Normalverteilung folgen.

7. ANOVA-Typ Statistik nach Werner ohne Normalverteilung

Solange es sich bei dem Versuchsaufbau um ein Ein-Stichproben-Design handelt, kann neben dem bereits hergeleiteten Ansatz mittels der empirischen Kovarianzmatrix auch der Weg über den Ansatz von Werner (2004) gegangen werden. Das Ziel bleibt weiterhin eine Teststatistik für die quadratische Form der Mittelwerte $Q_n^* = n\bar{\mathbf{Y}}'\bar{\mathbf{Y}}$ zu finden. Hierfür wird ebenfalls die, mit der Spur der Kovarianzmatrix skalierte quadratische Form Q_n^* aus (5.1) betrachtet, mit: $\frac{Q_n^*}{Sp(\mathbf{\Sigma})} \overset{d}{\sim} \chi_f^2/f$. Der Unterschied besteht darin, dass anstelle der Spur der empirischen Kovarianzmatrix der erwartungstreue Schätzer $B_0 = \frac{1}{n} \cdot \sum_{k=1}^n A_k$ zum Schätzen der Spur der Kovarianzmatrix verwendet wird.

Wobei sich B_0 ebenfalls als eine quadratische Form darstellen lässt und ohne Annahme der Normalverteilung Q_n^* und B_0 nicht unkorreliert sind (Argumentation analog zu Abschnitt 5.2). Daher wird die neue Zufallsvariable

$$D_n = \frac{n \cdot \bar{\mathbf{Y}}'\bar{\mathbf{Y}}}{B_0} = \frac{Q_n^*}{B_0} \quad (7.1)$$

definiert, dessen Verteilung unbekannt ist. Da aber $B_0 \xrightarrow{p} Sp(\mathbf{\Sigma})$ konvergiert, kann die unbekannte Verteilung des Quotienten D_n analog zu Abschnitt 5.2 mit Hilfe einer Box-Approximation angenähert werden.

Dies geschieht, indem die ersten beiden Momente von D_n mit denen einer um g_2 gestreckten $\chi_{f_2}^2/f_2$ -Verteilung gleichgesetzt werden. Um die Momente des Quotienten D_n anzunähern, wird ebenfalls die in Lemma 5.3.1 vorgestellte Taylor-Approximation verwendet. Doch vorab werden die nötigen Eigenschaften des Schätzers B_0 unter den in (2.1) gemachten Modellannahmen hergeleitet.

7.1. Der Schätzer B_0

Lemma 7.1.1 Für die Zufallsvektoren gelte $\mathbf{X}_k = (X_{k1}, \dots, X_{kd})' = \mathbf{\Gamma}\mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}$, $k = 1, \dots, n$, wie in (2.1) mit $\text{Cov}(\mathbf{\Gamma}\mathbf{Z}_k) = \mathbf{\Gamma}\mathbf{\Gamma}' = \mathbf{S}$ und $E(Z_{ks}^4) \leq \gamma < \infty \forall k = 1, \dots, n, s = 1, \dots, d$. Sei $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k + \mathbf{T}\boldsymbol{\mu}$ definiert wie in Proposition 2.3.1 mit $\text{Cov}(\mathbf{Y}_k) = \mathbf{T}\mathbf{S}\mathbf{T}' = \boldsymbol{\Sigma}$. Weiterhin sei $A_k = \mathbf{X}_k' \mathbf{T}\mathbf{X}_k = \mathbf{Y}_k' \mathbf{Y}_k$, $k = 1, \dots, n$ die quadratische Form der Zufallsvektoren (Siehe Definition 6.1.1) und $B_0 = \frac{1}{n} \sum_{k=1}^n A_k$.

Dann gilt unter $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$: B_0

1. ist ein erwartungstreuer Schätzer für $Sp(\boldsymbol{\Sigma})$.
2. ist konsistent im Sinne von Definition A.3.2.
3. ist dimensionsstabil im Sinne von Definition A.3.4.

Beweis:

Unter H_0 : $\mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ folgt aus Proposition 2.3.1:

$\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k$ sowie $E_{H_0}(\mathbf{Y}_k) = \mathbf{0}$.

Damit ergeben sich die folgenden Resultate:

1.

Die Erwartungstreue wird mit Hilfe von Lemma 6.1.3 (Momente I) nachgewiesen:

$$E_{H_0}(B_0) = E\left(\frac{1}{n} \sum_{k=1}^n A_k\right) = \frac{1}{n} \sum_{k=1}^n E(A_k) = Sp(\boldsymbol{\Sigma}).$$

2. und 3.

Zum Nachweis der Dimensionsstabilität und Konsistenz wird eine Abschätzung für die Varianz des Schätzers, welche mit Hilfe von Lemma 6.1.8 (3) zu bestimmen ist, benötigt.

$$\begin{aligned} \text{Var}(B_0) &= \text{Var}\left(\frac{1}{n} \sum_{k=1}^n A_k\right) = \frac{1}{n} \text{Var}(A_1) \\ &\leq \frac{1}{n} [\gamma \cdot Sp(\boldsymbol{\Sigma}^2)] \end{aligned}$$

Nun lässt sich auch die Dimensionsstabilität des Schätzers B_0 nachweisen, wobei dafür Lemma 6.1.3 (5) benutzt wird:

$$\begin{aligned}
 \text{Var} \left(\frac{B_0}{Sp(\boldsymbol{\Sigma})} \right) &= \frac{\text{Var}(B_0)}{[Sp(\boldsymbol{\Sigma})]^2} \\
 &\leq \frac{\frac{1}{n} [\gamma \cdot Sp(\boldsymbol{\Sigma}^2)]}{[Sp(\boldsymbol{\Sigma})]^2} \leq \frac{\frac{1}{n} [\gamma \cdot [Sp(\boldsymbol{\Sigma})]^2]}{[Sp(\boldsymbol{\Sigma})]^2} \\
 &= \frac{1}{n} \cdot \gamma = C_0(n).
 \end{aligned}$$

$C_0(n)$ ist somit unabhängig von der Wahl der Dimension d und für $n \rightarrow \infty$ geht $C_0(n) \rightarrow 0$. Damit ist die Dimensionsstabilität im Sinne von Definition A.3.4 erfüllt. Da $C_0(n)$ eine Nullfolge ist, folgt auch die Konsistenz im Sinne von Definition A.3.2.

□

Da auch der Schätzer B_0 alle geforderten Eigenschaften erfüllt, kann jetzt analog zur Herleitung in Kapitel 5 vorgegangen werden.

7.2. Momente der quadratischen Formen

Es werden zuerst die für die Taylor-Approximation der Momente benötigten Erwartungswerte, Varianzen und Kovarianzen mit Hilfe der Sätze 5.3.2 und 5.3.3 bestimmt. Dafür ist es von Vorteil, den Schätzer B_0 ebenfalls als quadratische Form darzustellen:

$$B_0 = \frac{1}{n} \sum_{k=1}^n A_k = \frac{1}{n} \cdot \mathbf{Y}' (\mathbf{I}_n \otimes \mathbf{I}_d) \mathbf{Y} = \frac{1}{n} \cdot \mathbf{Y}' (\mathbf{C}_n \otimes \mathbf{I}_d) \mathbf{Y} = \mathbf{Y}' \mathbf{C}^* \mathbf{Y}$$

mit $\mathbf{C}^* = \frac{1}{n} \cdot (\mathbf{C}_n \otimes \mathbf{I}_d) = \frac{1}{n} \cdot (\mathbf{I}_n \otimes \mathbf{I}_d)$ und

$$Q_n^* = n \cdot \bar{\mathbf{Y}}' \bar{\mathbf{Y}} = \mathbf{Y}' \left(\frac{1}{n} \mathbf{J}_n \otimes \mathbf{I}_d \right) \mathbf{Y} = \mathbf{Y}' \mathbf{A} \mathbf{Y}$$

mit $\mathbf{A} = (\mathbf{A}_n \otimes \mathbf{I}_d) = \left(\frac{1}{n} \mathbf{J}_n \otimes \mathbf{I}_d \right)$ wie bereits in Abschnitt 5.3 definiert.

Es fehlen lediglich die Varianz des Schätzers B_0 und die Kovarianz von B_0 und Q_n^* . Da die Voraussetzungen für die Sätze 5.3.2 und 5.3.3 erfüllt sind (siehe Abschnitt 5.4), lassen sich Varianz und Kovarianz wie folgt darstellen:

$$\begin{aligned}
 \text{Var}(B_0) &= \text{Var}(\mathbf{Y}' \mathbf{C}^* \mathbf{Y}) = \frac{1}{n^2} \cdot \text{Var}(\mathbf{Y}' (\mathbf{I}_n \otimes \mathbf{I}_d) \mathbf{Y}) \\
 &= \frac{1}{n^2} \left[\left(\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 - 2Sp(\boldsymbol{\Sigma}^2) \right) \cdot n + 2Sp(\boldsymbol{\Sigma}^2) \cdot n \right] \\
 &= \frac{1}{n} \cdot \left(\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 - 2Sp(\boldsymbol{\Sigma}^2) \right) + 2Sp(\boldsymbol{\Sigma}^2) \cdot \frac{1}{n}
 \end{aligned}$$

$$= \frac{1}{n} \cdot \left[\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 \right] \text{ und}$$

$$\begin{aligned} Cov(Q_n^*, B_0) &= Cov(\mathbf{Y}' \mathbf{A} \mathbf{Y}, \mathbf{Y}' \mathbf{C}^* \mathbf{Y}) = \frac{1}{n} \cdot Cov(\mathbf{Y}' \left(\frac{1}{n} \mathbf{J}_n \otimes \mathbf{I}_d \right) \mathbf{Y}, \mathbf{Y}' (\mathbf{I}_n \otimes \mathbf{I}_d) \mathbf{Y}) \\ &= \frac{1}{n} \cdot \left[\left(\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 - 2 \cdot Sp(\boldsymbol{\Sigma}^2) \right) \cdot 1 + 2 \cdot Sp(\boldsymbol{\Sigma}^2) \cdot 1 \right] \\ &= \frac{1}{n} \cdot \left[\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 \right] = Var(B_0). \end{aligned}$$

7.2.1. Taylor-Approximation

Wie angekündigt, werden jetzt die ersten beiden Momente des Quotienten D_n mit Hilfe der Taylor-Approximation (Lemma 5.3.1) angenähert. Die dafür benötigten Erwartungswerte, Varianzen und Kovarianzen der quadratischen Formen Q_n^* und B_0 werden vorab aufgelistet.

$$\begin{aligned} E(Q_n^*) &= Sp(\boldsymbol{\Sigma}) \\ E(B_0) &= Sp(\boldsymbol{\Sigma}) \\ Var(Q_n^*) &= \frac{1}{n} \cdot \left(\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 - 2Sp(\boldsymbol{\Sigma}^2) \right) + 2Sp(\boldsymbol{\Sigma}^2) \\ Var(B_0) &= \frac{1}{n} \cdot \left[\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 \right] \\ Cov(Q_n^*, B_0) &= \frac{1}{n} \cdot \left[\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 \right] \end{aligned}$$

Nun ist auch die Taylor-Approximation ohne weiteres durchführbar und zur Vereinfachung der Schreibweise werden wieder folgende Bezeichnungen verwendet:

$$\beta_1 = [Sp(\boldsymbol{\Sigma})]^2 \text{ und } \beta_2 = Sp(\boldsymbol{\Sigma}^2).$$

Approximation des Erwartungswertes und der Varianz von $D_n = \frac{Q_n^*}{B_0}$:

$$\begin{aligned} E\left(\frac{Q_n^*}{B_0}\right) &\doteq \frac{E(Q_n^*)}{E(B_0)} \left[1 + \frac{Var(B_0)}{[E(B_0)]^2} - \frac{Cov(Q_n^*, B_0)}{E(Q_n^*) E(B_0)} \right] \\ &= \frac{Sp(\boldsymbol{\Sigma})}{Sp(\boldsymbol{\Sigma})} \left[1 + \frac{\frac{1}{n} [\tau_4 - \beta_1]}{\beta_1} - \frac{\frac{1}{n} [\tau_4 - \beta_1]}{\beta_1} \right] \\ &= 1, \end{aligned}$$

$$\begin{aligned}
 \text{Var} \left(\frac{Q_n^*}{B_0} \right) &\doteq \frac{[E(Q_n^*)]^2}{[E(B_0)]^2} \left[\frac{\text{Var}(Q_n^*)}{[E(Q_n^*)]^2} + \frac{\text{Var}(B_0)}{[E(B_0)]^2} - 2 \frac{\text{Cov}(Q_n^*, B_0)}{E(Q_n^*) E(B_0)} \right] \\
 &= \frac{\beta_1}{\beta_1} \left[\frac{\frac{1}{n}(\tau_4 - \beta_1 - 2\beta_2) + 2\beta_2}{\beta_1} + \frac{\frac{1}{n}[\tau_4 - \beta_1]}{\beta_1} - 2 \frac{\frac{1}{n}[\tau_4 - \beta_1]}{\beta_1} \right] \\
 &= \frac{2\beta_2 \cdot \frac{(n-1)}{n}}{\beta_1} = \frac{2Sp(\Sigma^2) \cdot \frac{(n-1)}{n}}{[Sp(\Sigma)]^2}.
 \end{aligned}$$

7.2.2. Box-Approximation

Mit den Ergebnissen aus der Taylor-Approximation ist es möglich die Verteilung des Quotienten $D_n = \frac{Q_n^*}{B_0}$ mittels einer Box-Approximation anzunähern.

Dies geschieht analog zur Box-Approximation in Abschnitt 5.4.2.

Sei dafür:

$$U_2 \sim \chi_{f_2}^2 / f_2$$

mit

$$E(D_n) = E \left(\frac{Q_n^*}{B_0} \right) \doteq 1 = E(g_2 U_2) = g_2,$$

$$\text{Var}(D_n) = \text{Var} \left(\frac{Q_n^*}{B_0} \right) \doteq \frac{2 \cdot Sp(\Sigma^2) \cdot \frac{(n-1)}{n}}{[Sp(\Sigma)]^2} = \text{Var}(g_2 U_2) = 2 \cdot \frac{g_2^2}{f_2}.$$

Der Streckungsparameter $g_2 = 1$ kann direkt in die zweite Gleichung eingesetzt werden.

Es folgt die Berechnung des Freiheitsgrades f_2 :

$$\begin{aligned}
 f_2 &= \frac{2}{\text{Var}(g_2 U_2)} = \frac{n}{(n-1)} \cdot \frac{[Sp(\Sigma)]^2}{Sp(\Sigma^2)} \\
 &= \frac{n}{(n-1)} \cdot f,
 \end{aligned}$$

wobei $f = \frac{[Sp(\Sigma)]^2}{Sp(\Sigma^2)}$ der Freiheitsgrad aus der Boxapproximation (5.1) mit bekannter Kovarianzmatrix ist. Für $n \rightarrow \infty$ geht der Freiheitsgrad $f_2 \rightarrow f$.

7. ANOVA-Typ Statistik nach Werner ohne Normalverteilung

Das heißt, es ist möglich die Statistik D_n durch eine mit g_2 gestreckte $\chi_{f_2}^2/f_2$ -Verteilung zu approximieren, wobei für $n \rightarrow \infty$ D_n approximativ einer χ_f^2/f -Verteilung folgt (siehe (5.1)).

$$\frac{D_n}{g_2} = \frac{Q_n^*}{B_0} \dot{\sim} \chi_{f_2}^2/f_2 \quad f_2 = \frac{n}{(n-1)} \cdot \frac{[Sp(\boldsymbol{\Sigma})]^2}{Sp(\boldsymbol{\Sigma}^2)} \quad (7.2)$$

Es besteht auch hier das Problem, dass f_2 die unbekannt Parameter $[Sp(\boldsymbol{\Sigma})]^2$ und $Sp(\boldsymbol{\Sigma}^2)$ enthält, welche im Folgenden geschätzt werden müssen.

7.3. Die ANOVA-Typ Statistik nach Werner

Wie in Kapitel 6 gezeigt, sind die Schätzer B_1 und B_2 erwartungstreue, konsistente und dimensionsstabile Schätzer für $[Sp(\Sigma)]^2$ und $Sp(\Sigma^2)$. Durch Einsetzen der Schätzer B_1 und B_2 in den Freiheitsgrad f_2 ergibt sich der Schätzer:

$$\hat{f}_2 = \frac{n}{(n-1)} \cdot \frac{B_1}{B_2}.$$

Durch Verwendung des eben definierten Schätzers \hat{f}_2 in (7.2) ergibt sich ebenfalls eine Teststatistik:

$$A_n^{ATS-W} = D_n = \frac{Q_n^*}{B_0} \dot{\sim} \chi_{\hat{f}_2}^2 / \hat{f}_2, \quad \hat{f}_2 = \frac{n}{(n-1)} \cdot \frac{B_1}{B_2}. \quad (7.3)$$

Diese Statistik ist aber gerade die in (4.5) vorgestellte ATS-Werner Statistik, nur ohne explizite Annahme der Normalverteilung. Daraus folgt, dass die in Werner (2004) benötigte Annahme der Normalverteilung zum Nachweis der Dimensionsstabilität des Schätzers B_2 und zur Berechnung der Varianz zu streng gewählt wurde.

Im Einstichprobenfall konnten somit zwei Statistiken unter den Modellannahmen (2.1) auf Seite 5 hergeleitet werden, welche im nachfolgenden Kapitel miteinander verglichen werden sollen.

8. Simulationen der Statistiken

8.1. Simulationstechniken

In diesem Kapitel werden Simulationen bezüglich der hergeleiteten Teststatistiken durchgeführt, um das Verhalten dieser Tests unter verschiedenen konstruierten Bedingungen zu untersuchen. Zu diesem Zweck werden zu unterschiedlichen Modellkonstellationen Niveau- und Powersimulationen der Teststatistiken durchgeführt.

In den folgenden Abschnitten wird zuerst die Vorgehensweise der einzelnen Simulationen beschrieben, sowie die Erzeugung der Zufallszahlen erläutert und abschließend werden die Ergebnisse der Niveau- und Powersimulationen vorgestellt.

8.1.1. Niveau und Power

Es werden Zufallszahlen mit vorgegebenen Modellparametern erzeugt und anschließend die Teststatistiken darauf angewendet. Dabei ist es das Ziel, Rückschlüsse auf das Verhalten der Statistiken unter bestimmten Bedingungen zu erhalten. Ist die Anzahl der Wiederholungen groß genug, so lassen sich aus diesen Simulationen die Charakteristika der Tests und eventuelle Problembereiche herauslesen.

Bei Niveausimulationen werden die Zufallszahlen unter Hypothese erzeugt und zu einem vorgegebenen Niveau α getestet. Es wird dafür in jedem Durchlauf der p-Wert der Teststatistik bestimmt und mit α verglichen. Die relative Häufigkeit der p-Werte \hat{p} mit $\hat{p} < \alpha$ ist ein guter Schätzer für das wahre Niveau des Tests, wobei für eine *gute* Statistik das empirische Niveau möglichst nahe an dem vorgegebenen Niveau α liegen sollte. Mit wachsender Anzahl an Simulationsdurchläufen erhöht sich auch die Genauigkeit des Schätzers. In dieser Arbeit wurden alle Niveausimulationen zu $\alpha = 0,05$ mit $n_{sim} = 10.000$ Wiederholungen durchgeführt.

Um die Güte der jeweiligen Teststatistik besser darstellen zu können, wird für den Niveauschätzer \hat{p} ein Zufallsstreifen berechnet. Dieser wird mit Hilfe des Zentralen Grenzwertsatzes (siehe z.B. Dehling (2004), S. 213-226) hergeleitet, da jede Testentscheidung als Bernoulli-verteilte Zufallsvariable mit Erfolgswahrscheinlichkeit α aufgefasst werden kann und somit \hat{p} die relative Häufigkeit von $n_{sim} = 10.000$ Bernoulli-verteilten Zufallsvariablen ist. Es folgt die allgemeine Darstellung des Zufallsstreifens:

$$ZS_{\hat{p}} = \left[\alpha - \frac{1}{\sqrt{n_{sim}}} \cdot \sqrt{\alpha \cdot (1 - \alpha)} \cdot u_{1 - \frac{\alpha_z}{2}} ; \alpha + \frac{1}{\sqrt{n_{sim}}} \cdot \sqrt{\alpha \cdot (1 - \alpha)} \cdot u_{1 - \frac{\alpha_z}{2}} \right],$$

wobei $u_{1 - \frac{\alpha_z}{2}}$ das $1 - \frac{\alpha_z}{2}$ -Quantil der Standardnormalverteilung ist.

Damit ergibt sich für die durchgeführten Simulationen, zum Zufallsniveau $\alpha_z = 0,01$, folgender Zufallsstreifen:

$$ZS_{\hat{p}} = [0,04438 ; 0,05561].$$

Bei Powersimulationen wird das Verhalten der Tests unter Alternative (zur Hypothese $H_0: \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$) untersucht. Es ist dabei von Interesse, wie schnell ein Test eine vorgegebene Alternative aufdecken kann. Dafür werden analog zur Niveausimulation Zufallszahlen zu einem festen Niveau α erzeugt. Zusätzlich erhalten alle Zufallszahlen X_{ki} eine Verschiebung in der Form $X_{ki} + \mu_i$, wobei je nach Struktur der Powersimulation auch einige $\mu_i = 0$ sein dürfen. Diese Verschiebung wächst solange bis der Test die Verschiebung sicher aufdeckt. Die Power einer *guten* Teststatistik sollte also von α beginnend möglichst schnell gegen 1 gehen. Der verwendete Schätzer für die Power eines Tests ist, ähnlich zum Niveau, die relative Häufigkeit der p-Werte \hat{p} , nur diesmal mit $\hat{p} > \alpha$. In dieser Arbeit werden zwei Arten von Power untersucht, Ein-Punkt-Power und Trend-Power, wobei sich die Bezeichnungen auf die Strukturen der Alternative beziehen.

Im Falle der Ein-Punkt-Alternative, wird nur ein $\mu_i = \delta$ gesetzt und für $j \neq i = 1, \dots, d$ gilt $\mu_j = 0$, somit wird in allen Zufallsvektoren nur die i -te Komponente um δ verschoben. Dabei ist es für die Power der Teststatistiken unerheblich, welche Komponente verschoben wird.

Für die Trend-Alternative werden alle Zufallszahlen verschoben, wobei $\mu_i = \delta \cdot i/d$ gesetzt wird für $i = 1, \dots, d$ und das δ variiert. Es ergibt sich dadurch ein aufsteigender Trend.

Dabei ist zu beachten, dass die betrachteten Teststatistiken nicht für eine spezielle Art von Trend (z.B. aufsteigend oder absteigend) konzipiert wurden und somit die Reihenfolge der Verschiebungen μ_i , $i = 1, \dots, d$, keinen Einfluss auf die Power dieser Teststatistiken hat. Aus diesem Grund wird in den Folgenden Simulationen, repräsentativ für die Trend-Alternative, nur ein aufsteigender Trend betrachtet.

Mit diesen beiden Methoden kann das Verhalten der Tests auf lokale und globale Verschiebungen untersucht werden. In dieser Arbeit wurden alle Powersimulationen zu

$\alpha = 0,05$ mit $n_{sim} = 10.000$ Wiederholungen durchgeführt. Die Verschiebung δ wuchs dafür jeweils in kleinen Schritten, von 0 beginnend, bis die Alternativen sicher aufgedeckt wurden.

8.1.2. Kovarianzstrukturen

Im Rahmen der Simulationen werden weiterhin drei verschiedene Klassen von Kovarianzstrukturen untersucht, die entweder als spezielles Kovarianzmodell bei Messwiederholungen auftreten oder die allgemeine Robustheit der Tests aufzeigen sollen.

Compound Symmetry (CS): Diese Art von Kovarianzstruktur liegt vor, wenn zum Beispiel Zellkulturen eines Individuums unter verschiedenen Bedingungen gemessen werden. In diesem Fall haben alle Messungen dieselbe Varianz und alle Kovarianzen sind identisch (siehe Definition 2.2.1).

In solch einem Fall liegt eine $CS(\sigma^2, \tau)$ Struktur vor.

Autoregressiv (AR): Diese Art von Kovarianz tritt bei Messwiederholungen in Form von Zeitreihen auf, wobei davon ausgegangen wird, dass zeitlich näher zusammenliegende Messungen höher miteinander korreliert sind, als weiter voneinander entfernt liegende (siehe Definition 2.2.2).

In kurzer Schreibweise liegt dann eine $AR(\sigma^2, \rho)$ Struktur vor.

Unstrukturiert (UN): in diesem Fall folgt die Kovarianzmatrix keiner bestimmten Struktur.

8.1.3. Struktur der Zufallszahlen

Im Rahmen der Simulationen wurden in SAS Zufallsvariablen mit verschiedenen Verteilungen und unterschiedlichen Kovarianzstrukturen erzeugt. Pro d-dimensionalen Zufallsvektor wurden d unabhängige Zufallszahlen mit identischer Verteilung für den Messfehler erzeugt. Zusätzlich wurde auf jeden Zufallsvektor eine weitere unabhängig erzeugte Zufallszahl als interindividuelle Streuung addiert.

Damit ergibt sich für das Grund-Modell:

$$\mathbf{X} = \mathbf{Z} + E \cdot \mathbf{1}_d.$$

Die folgenden Verteilungen wurden für den Zufallsvektor \mathbf{X} erzeugt:

- Normalverteilung
- Log-Normalverteilung
- Gleichverteilung auf $[-1,1]$
- Exponentialverteilung mit $\lambda = 1$
- Bernoulli-Verteilung $p = 0,3$

Wird das Modell in seiner Grundform betrachtet, so liegt eine Compound Symmetry (CS) Struktur vor, da alle Varianzen sowie alle Kovarianzen identisch sind:

$$\text{Var}(\mathbf{X}) = \sigma^2 \mathbf{I}_d + \tau \mathbf{J}_d.$$

Um eine andere Kovarianzstruktur \mathbf{V} für den Messfehler zu erzeugen, wird der Ansatz aus (2.1) (Modell) verwendet, mit $\mathbf{A} = \mathbf{V}^{1/2}$:

$$\text{Var}(\mathbf{AZ}) = \mathbf{ASA}.$$

Dabei sollte $\mathbf{S} = \text{Var}(\mathbf{Z}) = \mathbf{I}_d$ gelten, ansonsten wird $\mathbf{A} = \frac{1}{\sigma} \cdot \mathbf{V}^{1/2}$ gewählt.

Um für den Messfehler eine autoregressive Kovarianzstruktur zu erzeugen, wird die Matrix \mathbf{A}_{AR} so gewählt, dass alle Kovarianzen die gewünschte Gestalt

$$\text{Cov}(X_{ki}, X_{kj}) = \rho^{|i-j|} \cdot \sigma^2$$

besitzen.

Das gewünschte Modell ergibt sich durch Multiplikation der Matrix \mathbf{A}_{AR} an den Zufallsvektor des Messfehlers \mathbf{Z} .

Analog verhält es sich mit der Erzeugung einer unstrukturierten Kovarianzmatrix (UN). Diese entsteht indem alle Einträge der Matrix \mathbf{A}_{UN} aus $\{-1,0,1\}$ -gleichverteilten Zufallszahlen bestehen, welche auf die Matrix $2 \cdot \mathbf{I}_d$ addiert werden. Das daraus resultierende Modell ergibt sich ebenfalls durch Multiplikation von \mathbf{A}_{UN} an den Vektor der Zufallszahlen des Messfehlers \mathbf{Z} .

Weiterhin ist für die korrekte Erzeugung der Zufallsvektoren eine Unterscheidung der verwendeten Verteilungen nötig.

1.)

Sind die zugrundegelegten Verteilungen des Zufallsvektors \mathbf{X} unendlich teilbar und ist die Darstellung des Erwartungswertes der Verteilung unabhängig von der Darstellung der Varianz und Kovarianz, so lässt sich \mathbf{X} als eine Summe von unabhängig identisch verteilten Zufallsvektoren darstellen und das Modell $\mathbf{X} = \mathbf{AZ} + E \cdot \mathbf{1}_d$ erzeugt die gewünschte Verteilung des Zufallsvektors \mathbf{X} mit entsprechender Kovarianzstruktur (z.B. Normalverteilung und Log-Normalverteilung).

Die Exponentialverteilung ist unendlich teilbar, da sie sich als eine Summe von zwei $\text{Gamma}(\frac{1}{2}, 1)$ -Verteilungen darstellen lässt, doch es ist für diese Verteilung nicht möglich, mit Hilfe des Modells $\mathbf{X} = \mathbf{AZ} + E \cdot \mathbf{1}_d$ eine autoregressive Kovarianzstruktur unter der Hypothese $H_0: \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ zu erzeugen.

2.)

Sind allerdings die simulierten Verteilungen nicht unendlich teilbar oder ist die Darstellung des Erwartungswertes abhängig von der Varianz und der Kovarianz (z.B. Gleichverteilung, Bernoulli-Verteilung und Exponentialverteilung), so werden für die Zufallsvektoren einige Transformationen benötigt bevor sie in den Simulationen verwendet werden können. Diese werden im Folgenden anhand der Gleichverteilung erläutert.

Um eine Kovarianzstruktur für einen gleichverteilten Zufallsvektor \mathbf{X} erzeugen zu können, werden zuerst normalverteilte Zufallsvektoren $\mathbf{AZ} + E \cdot \mathbf{1}_d$ mit compound symmetric oder autoregressiver Kovarianzstruktur erzeugt und anschließend durch Bestimmung der zugehörigen Quantile auf $(0,1)$ projiziert. Danach werden diese durch Anwenden der entsprechenden Umkehrfunktion in gleichverteilte Zufallsvektoren transformiert. Es ist dabei zu beachten, dass durch diese Transformationen nur die groben Eigenschaften der Kovarianzstrukturen erhalten bleiben und diese daher im Folgenden mit UNcs und UNar bezeichnet werden.

8.2. Niveau

Alle Simulationen wurden für verschiedene Statistiken im LD-F1 Design durchgeführt. Zum einen wurden die beiden in dieser Arbeit hergeleiteten Statistiken ANOVA-Typ Statistik nach Werner (ATS-Werner) (4.5) und die neue ANOVA-Typ Statistik (ATS-neu) (5.4) betrachtet, sowie zur Veranschaulichung die klassische ANOVA-Typ Statistik (ATS) (4.3) und die Box / Geisser-Greenhouse Statistik (GG) (4.4) angegeben. Die folgenden Grafiken zeigen die verschiedenen Simulationsergebnisse für die Teststatistiken, wobei für kleine Stichprobenumfänge $n = 10$ und zum Nachweis der Asymptotik $n = 30$ betrachtet werden. Die Anzahl der Messwiederholungen hingegen variiert von sehr klein $d = 3$ bis stark hochdimensional $d = 500$. Des Weiteren ist in allen Grafiken der Zufallsstreifen $ZS_{\hat{p}}$ durch zwei gestrichelte Linien gekennzeichnet. Es werden nur die wesentlichen Resultate der Simulationen dargestellt.

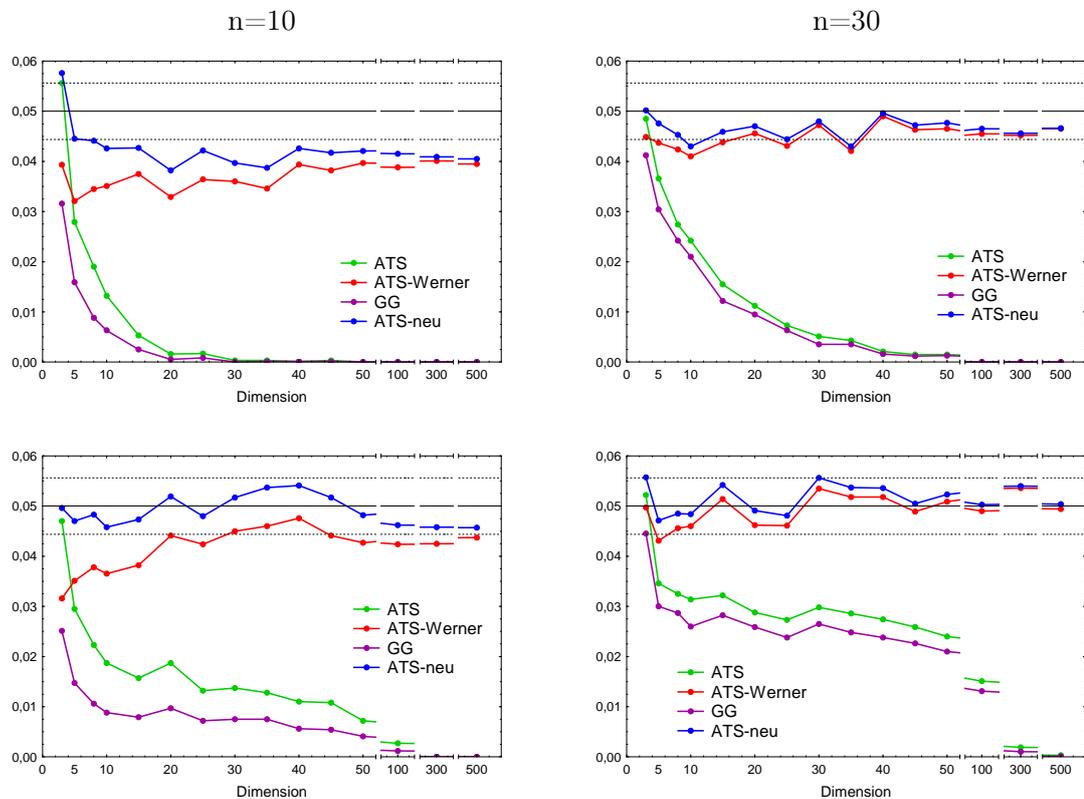


Abbildung 8.1.: Niveau: Exponentialverteilung, oben UNcs Struktur, unten UNar

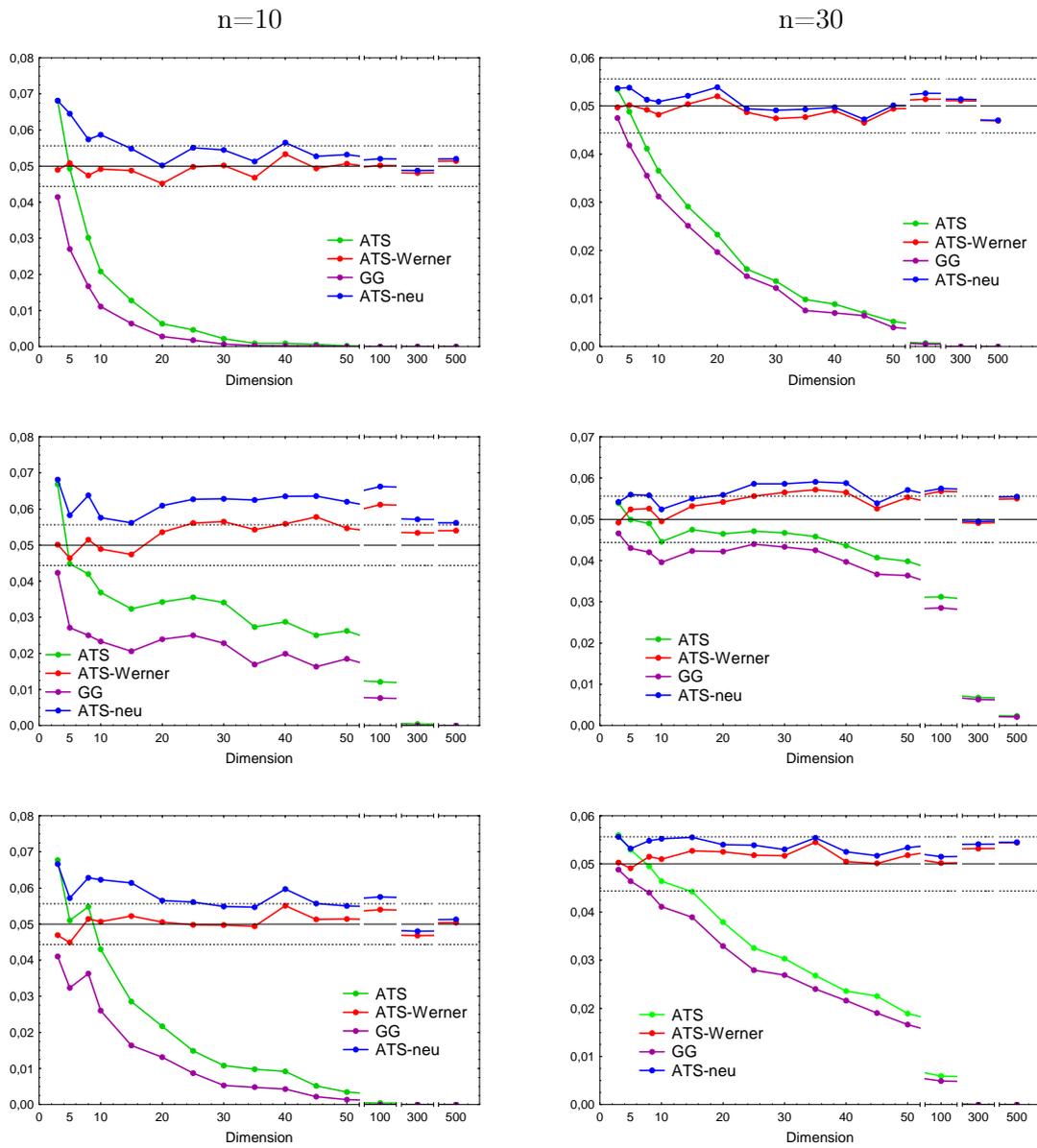


Abbildung 8.2.: Niveau: Normalverteilung, oben CS(2,1), mitte AR(1, 0,6), unten UN

8. Simulationen der Statistiken

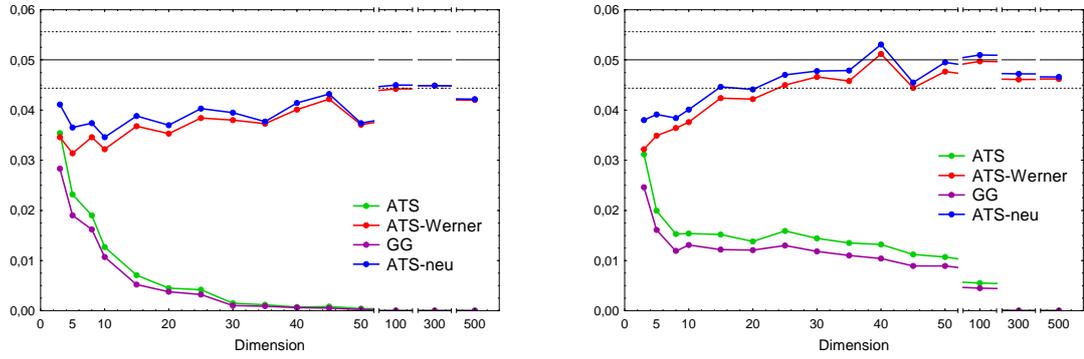


Abbildung 8.3.: Niveau: Log-Normalverteilung $n = 30$, links $CS(2,1)$, rechts $AR(1, 0,6)$

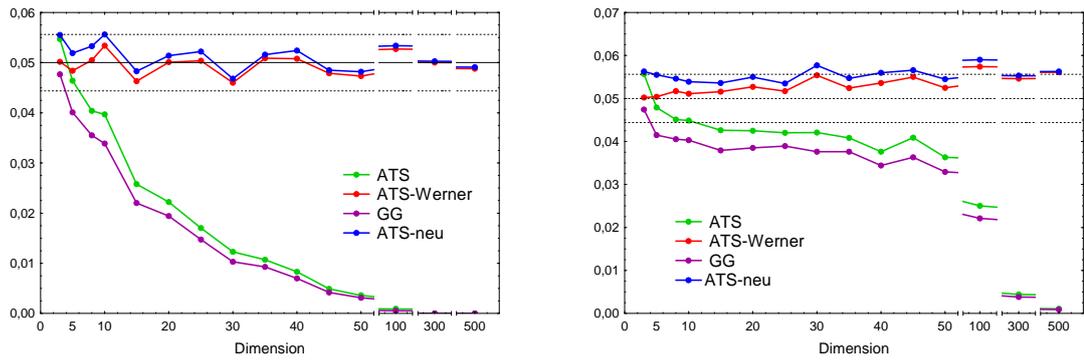
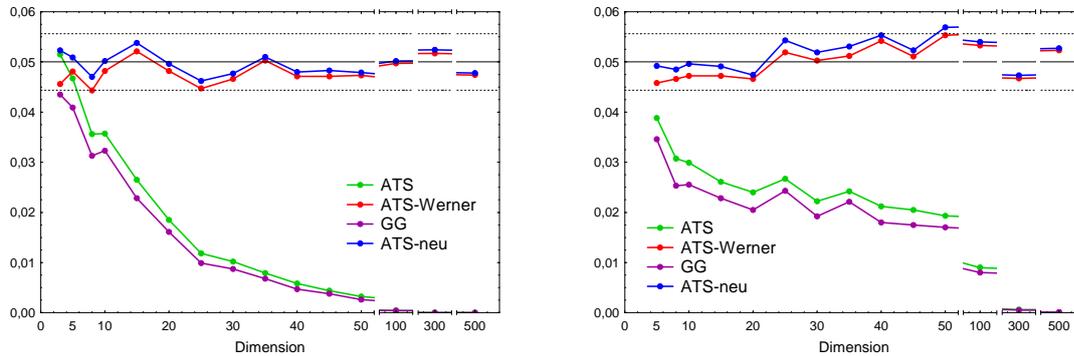


Abbildung 8.4.: Niveau: Gleichverteilung $n = 30$, links UNcs, rechts UNar

Abbildung 8.5.: Niveau: Bernulliverteilung $n=30$, links UNcs, rechts UNar

Wie aus den Niveausimulationen ersichtlich wird, hält die ATS-neu Statistik für $n = 30$ unter nahezu allen Verteilungen und Kovarianzstrukturen den Zufallsstreifen sehr gut ein, wobei die Güte der Approximation der Verteilung mit wachsender Anzahl an Messwiederholungen d immer weiter zunimmt. Für kleine Stichprobenumfänge ($n = 10$) und eine geringe Anzahl an Messwiederholungen d hingegen ist die Teststatistik leicht liberal, verbessert sich aber ebenfalls mit wachsendem d . Die zusätzlich betrachtete ATS-Werner Statistik hält unter vielen Verteilungen das Niveau schon für $n = 10$ sehr gut ein und die Güte der Teststatistik nimmt analog zur ATS-neu Statistik mit wachsender Dimension d zu. Es ist allerdings festzustellen, dass die ATS-Werner Statistik unter schiefen Verteilungen (z.B. Exponentialverteilung in Abbildung 8.1) mehr ins Konservative abrutscht als die neue ANOVA-Typ Statistik.

Die klassischen Verfahren ATS und GG, welche in den Abschnitten 4.2.1 und 4.2.2 vorgestellt wurden, halten im Repeated-Measures-Design unter keiner Verteilung das Niveau zufriedenstellend ein und sinken mit wachsender Dimension d immer weiter ins Konservative.

Die in Abbildung 8.5 dargestellten Simulationsergebnisse für Bernoulli-verteilte Zufallsvariablen wurden durchgeführt, um zu überprüfen in wie weit die Teststatistiken, welche metrische Messwerte voraussetzen, auf den Extremfall von dichotomen Daten reagieren. Es ist ersichtlich, dass beide Statistiken für $n = 30$ den Zufallsstreifen sehr gut einhalten und somit eine mögliche Erweiterung der Teststatistiken durch Ränge angestrebt werden könnte.

8.3. Power

In den Powersimulationen wurden analog zum Niveau die ATS-neu und ATS-Werner Statistik verwendet. Der Vergleichbarkeit halber waren die Stichprobenumfänge ebenfalls identisch bei $n = 10$ und $n = 30$, wohingegen die Anzahl der Messwiederholungen nur von $d = 3$ bis $d = 300$ betrachtet wurde. Da sich die beiden betrachteten Teststatistiken in allen Powersimulationen nahezu nicht voneinander unterschieden haben, wird in den folgenden Abbildungen nur die Power unter Normalverteilung für $n = 30$ und $d = 3, 10, 50, 100, 300$ dargestellt. Anhand dieser Darstellungen werden die verschiedenen Powersimulationen und Kovarianzstrukturen erläutert.

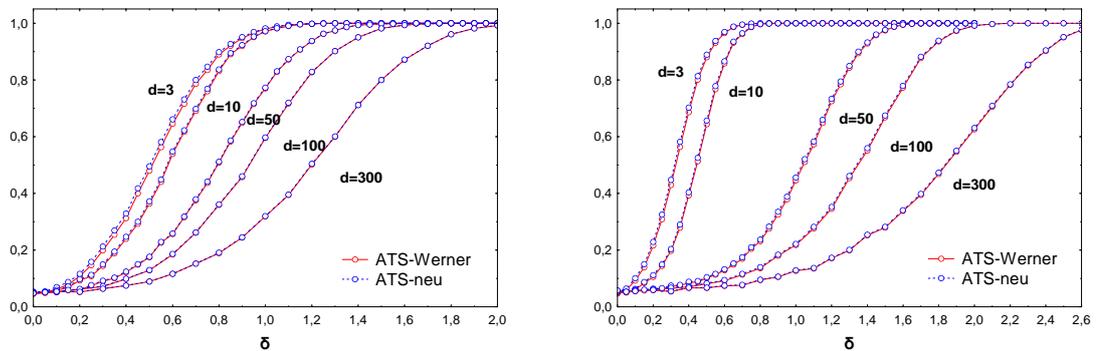


Abbildung 8.6.: Ein-Punkt-Power: Normalverteilung $n = 30$, links CS(2,1), rechts AR(1, 0,6)

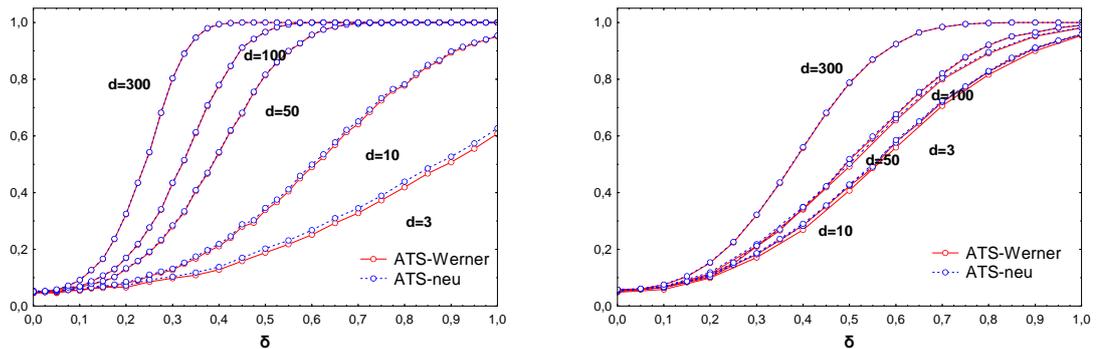


Abbildung 8.7.: Trend-Power: Normalverteilung $n = 30$, links CS(2,1), rechts AR(1, 0,6)

Wie aus den Abbildungen 8.6 und 8.7 ersichtlich, wird die Trend-Alternative für beide Statistiken mit wachsender Dimension d immer schneller aufgedeckt, wohingegen unter Ein-Punkt-Alternative die Power der Statistiken mit größer werdendem d immer langsamer wächst. Dies ist allerdings nicht ungewöhnlich, denn unter Ein-Punkt-Alternative wird bei einer sehr großen Anzahl von Messwiederholungen sprichwörtlich die Nadel im Heuhaufen gesucht. Weiterhin ist zu erkennen, dass für kleinere Dimensionen ($d < 50$) die Alternativen unter autoregressiver Kovarianzstruktur schneller aufgedeckt werden als unter Compound Symmetry. Entgegengesetzt verhält es sich für eine im Verhältnis zum Stichprobenumfang sehr große Anzahl an Messwiederholungen. Hier werden die Alternativen unter compound symmetric Kovarianzstruktur schneller aufgedeckt.

Im folgenden Abschnitt wird ein kurzer Überblick über die Handhabung der im Rahmen dieser Arbeit erweiterten Makros gegeben.

9. Software

Im Rahmen dieser Arbeit sind folgende Makros entstanden:

- NN-HD-F1
- NN-HD-F2
- NN-HD-F3

Diese Makros dienen dazu, ein- bis dreifaktorielle Repeated-Measures-Versuchspläne mit SAS auszuwerten, wobei in jedem Makro die ANOVA-Typ Statistik nach Werner und die neue ANOVA-Typ Statistik implementiert sind. Exemplarisch für die drei Makros wird im Anhang A.1 der Quelltext von NN_HD-F1 abgedruckt.

9.1. Einbinden und Aufrufen der Makros

Da die Verwendung aller Makros komplett identisch ist, wird hier die Handhabung anhand des NN-HD-F2 Makros erläutert.

Das Makro wird im SAS-Programmeditor mit dem Befehl

```
%INCLUDE 'C:\makro_NN_HD_F2'
```

eingebunden. Der auszuwertende Datensatz muss als SAS-Datei vorliegen wobei darauf zu achten ist, dass alle Messwerte als Vektor eingelesen werden müssen und nicht als Matrix. Somit ist es gegebenenfalls notwendig den Datensatz vor der Auswertung passend zu formatieren. Nachdem das Makro eingebunden wurde, wird es im SAS-Programmeditor mit dem Befehl

```
%NN_HD_F2( DATA = SAS-Datensatz,  
           VAR = SAS-Name der Zielvariable,  
           TIME1 = SAS-Name des ersten Zeitfaktors,  
           TIME2 = SAS-Name des zweiten Zeitfaktors,  
           SUBJEKT = SAS-Name der Individuen );
```

aufgerufen.

Die Zeitfaktoren stehen für die Struktur der verbundenen Messwiederholungen pro Individuum.

Dem Output sind dann die Werte für Statistiken, Freiheitsgrade und p-Werte zur Überprüfung der Hypothesen auf Effekte der Faktoren und der Wechselwirkung zu entnehmen, wobei alle Werte jeweils für die ATS-Werner und ATS-neu Statistik dargestellt werden.

Die Güte der in den Makros verwendeten Statistiken wurde bereits im Kapitel 8 (Simulationen) untersucht und belegt. Weiterhin wurden die Makros einer Reihe von Tests mittels konstruierter Datensätze unterzogen, um deren korrekte Funktionsweise zu überprüfen.

Mit Hilfe dieser Makros werden im folgenden Abschnitt die in Kapitel 4 vorgestellten Beispiele ausgewertet.

10. Auswertung der Beispiele

10.1. α -Amylase Studie

Der α -Amylase Versuch wurde mit Hilfe des NN_HD_F2 Makros und der neuen ANOVA-Typ Statistik (ATS-neu) ausgewertet. In den Klammern stehen die p-Werte der Werner-Statistik (ATS-Werner).

Tabelle 10.1.: Auswertung der α -Amylase Studie

| Untersuchung | Effekt | Statistik | df. | p-Wert |
|--------------|-------------------|-----------|--------|-----------------|
| Beide Tage | Tag | .01531 | 1.2363 | .94285 (.91185) |
| | Zeit | 5.3658 | 2.8133 | .00209 (.00700) |
| | Tag \times Zeit | 4.2168 | 2.4677 | .01258 (.02264) |
| Montag | Zeit | 1.7357 | 4.5015 | .2871 (.28994) |
| Donnerstag | Zeit | 8.1358 | 1.4751 | .00472 (.02374) |

Da eine Wechselwirkung zwischen den Tagen (*Montag*, *Donnerstag*) und dem Zeitverlauf über den Tag vorliegt, wird eine weitere Auswertung des Faktors Zeit getrennt nach den Tagen, mit Hilfe des NN_HD_F1 Makros und der neuen ANOVA-Typ Statistik, durchgeführt. Die nach Tagen getrennte Auswertung wurde zum $\alpha/2 = 2,5\%$ Niveau durchgeführt und die adjustierten p-Werte in Tabelle 10.1 dargestellt.

Abschließend ergibt sich, dass eine globale Schwankung der α -Amylase im Laufe eines Tages (Faktor: *Zeit*) signifikant festgestellt werden konnte, doch nach genauerer Untersuchung der Wechselwirkung wurde die Schwankung nur noch am Donnerstag, also in der Mitte der Woche, signifikant nachgewiesen. Daraus folgt, dass die Vermutung über den Einfluss der Tage auf den Tagesverlauf der α -Amylase statistisch bestätigt werden konnte, wohingegen eine spezifische tagesverlaufsbedingte Schwankung nur in der Mitte der Woche nachweisbar war. Weiterhin konnte kein globaler Effekt in der α -Amylase zwischen den Tagen festgestellt werden. Die Ergebnisse der Auswertung sind in Tabelle 10.1 dargestellt, wobei beide Statistiken zu denselben Testentscheidungen gekommen sind.

10.2. Cortisol-Konzentration im Blutplasma

Der Versuch über die Cortisol-Konzentration wurde mit Hilfe des NN_HD_F3 Makros und der neuen ANOVA-Typ Statistik (ATS-neu) ausgewertet. In den Klammern stehen die p-Werte der Werner-Statistik (ATS-Werner).

Tabelle 10.2.: Auswertung der Cortisol-Konzentration im Blutplasma

| Untersuchung | Effekt | Statistik | df. | p-Wert |
|-------------------------|-----------------|-----------|--------|-----------------|
| Vor/Nach Trainingspause | Zeit 1 | 6.0413 | 1.2803 | .01619 (.03577) |
| Placebo / m-CCP | Behandlung | 13.508 | 1.2803 | .00022 (.00842) |
| Messungen pro Tag | Zeit 2 | 2.5451 | 2.4768 | .08154 (.09069) |
| Wechselwirkungen | Zeit1 × Beh | .02241 | 1.2803 | .93424 (.89380) |
| | Zeit1 × Zeit2 | .78037 | 3.5967 | .54947 (.52467) |
| | Beh × Zeit2 | 4.1849 | 3.3213 | .00585 (.01361) |
| | Zeit1×Beh×Zeit2 | .47484 | 3.3180 | .73707 (.71564) |
| Placebo | Zeit 2 | 1.8618 | 3.7177 | .27180 (.27666) |
| m-CCP | Zeit 2 | .66338 | 4.8310 | .99999 (.99999) |

Aufgrund der Tatsache, dass eine statistische Wechselwirkung zwischen dem Faktor Behandlung (*Placebo*, *m-CCP*) und den Verlaufsmessungen pro Tag (Faktor: *Zeit 2*) besteht, werden diese Messungen nach den Behandlungsarten getrennt und mit dem NN_HD_F1 Makro und der ATS-neu Statistik ausgewertet. Analog zum ersten Beispiel, werden die nach Behandlung getrennten Auswertungen zum $\alpha/2 = 2,5\%$ Niveau durchgeführt und die adjustierten p-Werte in Tabelle 10.2 abgetragen.

Wie aus Tabelle 10.2 ersichtlich, ergibt sich eine signifikante Veränderung der Cortisol-Konzentration in Bezug auf die Trainingspause. Ebenfalls konnte ein signifikanter Einfluss des Mittels m-CCP gegenüber Placebo festgestellt werden. Eine signifikante Veränderung der Konzentration in den Verlaufsmessungen pro Tag konnte nicht festgestellt werden, allerdings hat die Behandlung einen Einfluss auf den Verlauf dieser Messungen. Somit konnte der vermutete Zusammenhang zwischen abrupten Trainingspausen und erhöhter Cortisol-Konzentration signifikant nachgewiesen werden.

Für weitere Informationen zu den beiden Datensätzen siehe Brunner (2002), S. 9-10.

11. Zusammenfassung und Ausblick

In dieser Arbeit wurden globale Testverfahren für den Einstichprobenfall im Repeated-Measures-Design, speziell für hochdimensionale Daten, ohne Annahme der Normalverteilung hergeleitet.

Nachdem die Probleme in den Herleitungen der bekannten Statistiken dargestellt wurden, welche auftreten, wenn die Annahme der Normalverteilung fallen gelassen wird oder das Versuchsdesign hochdimensional ist, konnten zwei dieser Verfahren weiter betrachtet werden. Um die Herleitungen dieser Verfahren auch ohne Normalverteilungsannahme durchführen zu können, wurde der zum Testen benötigte Quotient der quadratischen Formen als eine eigenständige Zufallsvariable angesehen. Im Zuge dieser Herleitungen entstanden die Sätze zur Darstellung der Varianz und Kovarianz von quadratischen Formen im Multivariaten. Des Weiteren wurden die gewünschten Eigenschaften der benötigten Schätzer (B_0, B_1, B_2) unter den gemachten Modellannahmen, speziell ohne Annahme der Normalverteilung, nachgewiesen. Auf diesem Weg konnte eine neue ANOVA-Typ Statistik hergeleitet und die ANOVA-Typ Statistik nach Werner verallgemeinert werden.

Darauf aufbauend wurden Makros für ein- bis dreifachstrukturierte Repeated-Measures-Designs in SAS weiterentwickelt, welche die hergeleiteten Teststatistiken enthalten. Die Funktionsweise der Makros wurde anhand von konstruierten Beispieldatensätzen überprüft und belegt. Die Teststatistiken selbst wurden in zahlreichen Simulationen auf ihr Verhalten unter verschiedenen Verteilungen und Modellannahmen untersucht.

Daraus ergab sich, dass für eine moderate Anzahl an Subjekten ($n=30$) beide Statistiken in fast allen Niveausimulationen und für eine beliebige Anzahl an Messwiederholungen innerhalb des Zufallsstreifens lagen. Weiterhin konnte beobachtet werden, dass sich die in den Teststatistiken verwendeten Approximationen der Verteilungen mit wachsender Anzahl an Messwiederholungen nicht verschlechtern haben. Die Teststatistiken hielten für eine große Anzahl an Messwiederholungen ($d \geq n$) den Zufallsstreifen sogar besser ein als für eine geringe Anzahl ($d < n$).

Damit ist es also gelungen, robuste Testverfahren für den Einstichprobenfall herzuleiten, welche auf eine breite Klasse von Verteilungen und für beliebig hochdimensionale Repeated-Measures-Designs anwendbar sind. Zusätzlich können diese Teststatistiken

bereits für eine sehr geringe Anzahl an Repeated-Measures ($d > 2$) verwendet werden.

Abschließend muss erwähnt werden, dass mit den in dieser Arbeit gemachten Modellannahmen keineswegs alle Verteilungen erreicht werden können, speziell nicht mit jeder beliebigen Kovarianzstruktur. Dennoch wird das Feld für die Anwendbarkeit dieser Statistiken erheblich über die Normalverteilung hinaus erweitert.

Nach der Herleitung der globalen Tests ohne Normalverteilung wäre der nächste Schritt, die Erweiterung der Statistiken auf den Mehrstichprobenfall, unter den gemachten Modellannahmen. Zu Beginn der Arbeit wurde diese Idee bereits angesprochen, wobei dafür der Weg analog zur neuen ANOVA-Typ Statistik (5.4) erstrebenswert wäre.

Ein weiterer Ansatz, wäre der Schritt hin zu multiplen Tests, denn bisher ist nur das Testen globaler Hypothesen möglich.

Ein anderer Ansatz zur Weiterentwicklung wäre der Versuch die bestehenden Modellannahmen noch weiter zu lockern und somit mehr Verteilungsklassen zu erschließen. Dabei muss untersucht werden, inwieweit Verallgemeinerungen möglich sind, ohne die, für Schätzer notwendige, Eigenschaft der Dimensionsstabilität zu verlieren. In diesem Zuge könnte auch versucht werden, die Approximationsgenauigkeit der Statistiken zu erhöhen, indem ein Restglied für die Taylor-Approximation bestimmt wird.

Ein weiterer Aspekt wäre der Übergang in die Nichtparametrik durch die Anwendung von Rängen auf die Messwiederholungen. Dies führt allerdings zu erheblichen Schwierigkeiten in den Herleitungen, da nach Anwendung der Ränge die Zufallsvektoren nicht mehr unabhängig sind.

A. Anhang

A.1. Makro

Exemplarisch für die drei verbesserten Makros ist hier der Übersichtlichkeit halber das Kürzeste abgedruckt:

```
/* **** */
/* *** Makro NN_HD_F1 *** */
/* *** überarbeitet von Hans-Joachim Helms, April 2010 *** */
/* *** */
/* **** */
%macro NN_HD_F1(data =,
                var =,
                TIME1 =,
                SUBJECT =);

proc sort data=&data ;
by &subject &time1 ;
run ;

proc iml worksize=120 ;
RESET LINESIZE = 80 ;
%LET faktora = &TIME1 ;
%LET faktorb = &SUBJECT ;

/* Funktion, welche die Box-Approximation durchführt */
/* ****matrix a kreuz n**** */
START box_neu(n,anzahl,z,C,chi2_helms,chi2_werner,fneu_werner,
fneu_helms,p_alt,p_neu,invers,d);
***** Zentrieren *****;
x = z' ;
T = C' * GINV(C * C') * C ;
xquer = x[+,]/n ;
*****;
help = j(n,d,0) ;
do i = 1 to n ;
do j = 1 to d ;
help[i,j] = x[i,j]-xquer[j] ;
```

```

end ;
end ;
/*****
/**** Stichproben-Kovarianzmatrix *****/
stich = (1/(n-1))*help'*help ;
****quadrat- und bilinearformen *****/
m = x * T * x' ;
bb_vektor = vecdiag(m) ;
bb_quad = bb_vektor * bb_vektor' ;
bb_test = bb_quad # invers ;
bil_test = m # invers ;
**** Statistiken *****/
zaehler_chi2 = n * xquer * T * xquer' ;
vau = ( x' * x ) / n ;
nenner_chi2_werner = trace ( T*vau ) ;
nenner_chi2_helms = trace ( T*stich ) ;
chi2_werner = zaehler_chi2/nenner_chi2_werner ;
chi2_helms = zaehler_chi2/nenner_chi2_helms ;
**** Freiheitsgrade *****/
**** Schätzer B1 *****/
b1 = 1/(n*(n-1)) * sum(bb_test) ;
aklquad = ssq(bil_test) ;
**** Schätzer B2 *****/
b2 = 1/(n*(n-1)) * aklquad ;
*** Freiheitsgrad für Werner *****/
fneu_werner = b1/b2/(1-1/n) ;
*** Freiheitsgrad für Helms *****/
fscho = (b1 + 2*b2/(n-1))*(b1 + 2*b2/(n-1)) ;
fschu = (b1 * b2 * (n / (n-1))) ;
fneu_helms = fscho / fschu ;
**** Streckungsparameter für Helms *****/
gneu_helms = (1 + (2*b2)/((n-1)*b1)) ;
**** p-Werte *****/
if chi2_werner > 0 & fneu_werner > 0 then
p_alt = 1-probchi(chi2_werner*fneu_werner, fneu_werner) ;
else p_alt = . ;
if chi2_helms > 0 & fneu_helms > 0 then
p_neu= 1-probchi((chi2_helms*fneu_helms)/gneu_helms, fneu_helms) ;
else p_neu = . ;
FINISH ;
reset nolog ;

```

```

/***** Daten einlesen *****/
USE &data ;
READ ALL VAR{&var} INTO werte ;
READ ALL VAR{&subject} INTO pat_ ;
READ ALL VAR{&time1} INTO t1_ ;
CLOSE &data ;

lev_a = unique(t1_) ; /* Stufen des Faktors A */
lev_b = unique(pat_) ; /* Stufen des Faktors B */
a = ncol(lev_a) ; /* Anzahl der Stufen von A */
n = ncol(lev_b) ; /* Anzahl der subjects */
d = a ; /* Anzahl der aller Stufen */
NN = a*n ; /* Anzahl der Daten */
print a n ;
/*****
/**** i ist der Laufindex des Faktors A i=1,...,a ****
/**** k ist der Laufindex der Subjects k=1,...,n ****
/****
/**** Definitionen ****
ea = j(a,1,1) ;
en = j(n,1,1) ;
ja = j(a,a,1) ;
jn = j(n,n,1) ;
pa = i(a) - ja/a ;
pn = i(n) - jn/n ;
invers = jn - i(n) ;
/**** Hypothesen-Matrizen ****
ma = pa ;
/****
***einlesen der MW in eine transponierte n x a Matrix***;
wert_matrix = (shape (werte,n,a))' ;
/****
Die BOX-Approximation: Ergebnisse in 'box'
*****/
box = j(1,6,0) ;
RUN box_neu(n,a,wert_matrix,ma,QF_Helms,QF_Werner,
DF_alt,DF_neu,p_alt,p_neu,invers,d) ; /* Effekt A*/
box[1,1]=QF_Werner; box[1,2]=DF_alt; box[1,3]=p_alt;
box[1,4]=QF_Helms; box[1,5]=DF_neu; box[1,6]=p_neu;

```

```

/*****
/***   FUNKTION                               ***/
/***   Ausdruck der Class Level Information (CLI) ***/
/*****/
start O_CLI (a,n,NN)                               ;
reset center                                       ;
class = {Time1 &faktora}                           ;
levels = a                                         ;
print /                                             ;
print 'New LD_F1 — subjects x Time1',             ;
      'Time1: fixed , subjects: random'           ;
print 'SAS-datafile-name:  ' "&data" ,             ;
      'Response variable:  ' "&var"                 ;
print 'Class Level Information'                   ;
print class levels                                 ;
reset noname                                       ;
print 'Total number of observations ' NN ' ',      ;
      'Total number of subjects   ' n ' ',        ;
finish                                             ; /* Ende von O_CLI */
/*****/
/* nu_char FUNKTION                               */
/* Diese Funktion schreibt numerische Werte in   */
/* characters um und trimmt diese durch max(...) */
/* bzw. trimmt diese , falls sie schon characters sind.*/
/*****/
start nu_char(v,v_neu)                             ;
if type (v) ='N' then v_neu =                      ;
char(v,max(int (log10(max(abs(v)))))+1))           ;
else v_neu = trim(v)                               ;
finish                                             ; /* Ende der Funktion nu_char */
/*****/
/**** Ausgabe der Quadratformen und der p_values ****/
/****/
start test_out(box)                                 ;
print 'Chi-Quadrat-Approximation'                 ;
source2 = {"Time1 " }                             ;
coll     = {"   QF_werner" "f_werner" "p-value(werner)" ;
           "QF_helms" "f_helms" "p-value(helms)"} ;
print 'Comparison of new ANOVA-type-statistic to old' ;
print box[r=source2][c=coll][format=6.5]         ;
finish                                             ;
/****/
/*Aufruf der Funktionen*****/
dataname = name(&data)                             ;

```

```
varname = name(&var)           ;
fa1_name = name(&faktora)      ;
fa2_name = name(&faktorb)      ;

run nu_char(lev_a, leva)       ;
run nu_char(lev_b, levb)       ;
run o_cli(a,n,NN)              ;
run test_out(box)              ;
quit                            ;
%mend NN_HD_F1                  ;
```

A.2. Beweise

Satz A.2.1 (Kovarianz von quadratischen Formen (Seite 25))

Die Zufallsvektoren $\mathbf{Y}_k = (Y_{k1}, \dots, Y_{kd})'$, $k = 1, \dots, n$, seien unabhängig und identisch verteilt und es bezeichne $\mathbf{Y} = (\mathbf{Y}'_1, \dots, \mathbf{Y}'_n)'$ den Vektor aller Zufallsvariablen mit $E_{H_0}(\mathbf{Y}) = \mathbf{0}$ und $Cov(\mathbf{Y}) = \mathbf{I}_n \otimes \Sigma$. Ferner sei $\tau_4 = E\left([\mathbf{Y}'_k \mathbf{Y}_k]^2\right)$, $k = 1, \dots, n$ und für alle $d < \infty$ soll $\tau_4 \leq \gamma_{\tau_4} < \infty$ gelten. Seien $\mathbf{A} = \mathbf{A}_n \otimes \mathbf{I}_d$ und $\mathbf{B} = \mathbf{B}_n \otimes \mathbf{I}_d$ symmetrische Matrizen mit $\mathbf{a}_n = \text{diag}\{\mathbf{A}_n\}$ und $\mathbf{b}_n = \text{diag}\{\mathbf{B}_n\}$.

Dann gilt:

$$Cov(\mathbf{Y}'\mathbf{A}\mathbf{Y}, \mathbf{Y}'\mathbf{B}\mathbf{Y}) = \left(\tau_4 - [Sp(\Sigma)]^2 - 2Sp(\Sigma^2)\right) \mathbf{a}'_n \mathbf{b}_n + 2Sp(\Sigma^2) Sp(\mathbf{A}_n \mathbf{B}_n).$$

Beweis:

Da $E_{H_0}(\mathbf{Y}) = \mathbf{0}$ und $Cov(\mathbf{Y}) = \mathbf{I}_n \otimes \Sigma$ gilt, stellt sich die Kovarianz zweier quadratischer Formen wie folgt dar:

$$Cov(\mathbf{Y}'\mathbf{A}\mathbf{Y}, \mathbf{Y}'\mathbf{B}\mathbf{Y}) = E(\mathbf{Y}'\mathbf{A}\mathbf{Y}\mathbf{Y}'\mathbf{B}\mathbf{Y}) - [E(\mathbf{Y}'\mathbf{A}\mathbf{Y})] \cdot [E(\mathbf{Y}'\mathbf{B}\mathbf{Y})],$$

wobei die Erwartungswerte der quadratischen Formen direkt über den Satz von Lancaster (siehe Satz A.3.6) bestimmt werden können.

$$E(\mathbf{Y}'\mathbf{A}\mathbf{Y}) = \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu} + Sp([\mathbf{A}_n \otimes \mathbf{I}_d] \cdot [\mathbf{I}_n \otimes \Sigma]) = Sp(\mathbf{A}_n) \cdot Sp(\Sigma)$$

$$E(\mathbf{Y}'\mathbf{B}\mathbf{Y}) = \boldsymbol{\mu}'\mathbf{B}\boldsymbol{\mu} + Sp([\mathbf{B}_n \otimes \mathbf{I}_d] \cdot [\mathbf{I}_n \otimes \Sigma]) = Sp(\mathbf{B}_n) \cdot Sp(\Sigma)$$

Daraus folgt für das Produkt der Erwartungswerte

$$[E(\mathbf{Y}'\mathbf{A}\mathbf{Y})] \cdot [E(\mathbf{Y}'\mathbf{B}\mathbf{Y})] = [Sp(\mathbf{A}_n) Sp(\mathbf{B}_n)] \cdot [Sp(\Sigma)]^2.$$

Somit muss nur noch die Darstellung von $E(\mathbf{Y}'\mathbf{A}\mathbf{Y}\mathbf{Y}'\mathbf{B}\mathbf{Y})$ gezeigt werden.

Da sich $\mathbf{Y}'\mathbf{A}\mathbf{Y}$ als eine Doppelsumme der \mathbf{Y}_k schreiben lässt, folgt daraus:

$$\begin{aligned} E(\mathbf{Y}'\mathbf{A}\mathbf{Y}\mathbf{Y}'\mathbf{B}\mathbf{Y}) &= E\left[\left(\sum_{i=1}^n\sum_{j=1}^na_{ij}\cdot\mathbf{Y}'_i\mathbf{Y}_j\right)\cdot\left(\sum_{r=1}^n\sum_{s=1}^nb_{rs}\cdot\mathbf{Y}'_r\mathbf{Y}_s\right)\right] \\ &= \sum_i^n\sum_j^n\sum_r^n\sum_s^n a_{ij}b_{rs}E(\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_r\mathbf{Y}_s), \end{aligned}$$

mit

$$E(\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_r\mathbf{Y}_s) = \begin{cases} \tau_4 & i = j = r = s \\ [Sp(\boldsymbol{\Sigma})]^2 & i = j \neq r = s \\ Sp(\boldsymbol{\Sigma}^2) & i = r \neq j = s \text{ und } i = s \neq j = r \\ 0 & \text{sonst.} \end{cases}$$

Die verschiedenen Fälle des Erwartungswertes werden im Folgenden bestimmt.

Da unter $H_0: \boldsymbol{\mu} = \mathbf{0}$ gilt, folgt daraus:

$$E(\mathbf{Y}'_1\mathbf{Y}_1) = \boldsymbol{\mu}'\mathbf{I}_d\boldsymbol{\mu} + Sp(\mathbf{I}_d \cdot \boldsymbol{\Sigma}) = Sp(\boldsymbol{\Sigma})$$

und im Einzelnen:

für $i = j = r = s$

$$E(\mathbf{Y}'_i\mathbf{Y}_i\mathbf{Y}'_i\mathbf{Y}_i) = E[(\mathbf{Y}'_i\mathbf{Y}_i)^2] = \tau_4$$

und mit $i = j \neq r = s$ folgt:

$$E(\mathbf{Y}'_i\mathbf{Y}_i\mathbf{Y}'_r\mathbf{Y}_r) = E(\mathbf{Y}'_i\mathbf{Y}_i) \cdot E(\mathbf{Y}'_r\mathbf{Y}_r) = [Sp(\boldsymbol{\Sigma})]^2.$$

Ist $i = s \neq j = r$, dann folgt:

$$\begin{aligned} E(\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_j\mathbf{Y}_i) &= E[Sp(\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_j\mathbf{Y}_i)] = E[Sp(\mathbf{Y}_i\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_j)] \\ &= Sp[E(\mathbf{Y}_i\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_j)] = Sp[E(\mathbf{Y}_i\mathbf{Y}'_i) \cdot E(\mathbf{Y}_j\mathbf{Y}'_j)] \\ &= Sp(\boldsymbol{\Sigma}^2). \end{aligned}$$

Analog gilt dies auch für $i = r \neq j = s$, da

$$\begin{aligned} E(\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_i\mathbf{Y}_j) &= E(\mathbf{Y}'_i\mathbf{Y}_j(\mathbf{Y}'_i\mathbf{Y}_j)) = E(\mathbf{Y}'_i\mathbf{Y}_j(\mathbf{Y}'_j\mathbf{Y}_i)) \\ &= E(\mathbf{Y}'_i\mathbf{Y}_j\mathbf{Y}'_j\mathbf{Y}_i) = Sp(\boldsymbol{\Sigma}^2). \end{aligned}$$

Nun ist es möglich, die oben definierte Summe in die jeweiligen Untersummen aufzuspalten, wobei sich alle a_{ij} , b_{rs} auf die Matrizen \mathbf{A}_n , \mathbf{B}_n und nicht auf \mathbf{A} , \mathbf{B} beziehen.

$$E[(\mathbf{Y}'\mathbf{A}\mathbf{Y})(\mathbf{Y}'\mathbf{B}\mathbf{Y})] = \tau_4 \sum_{i=1}^n a_{ii}b_{ii} + [Sp(\boldsymbol{\Sigma})]^2 \sum_{i \neq j} a_{ii}b_{jj} + Sp(\boldsymbol{\Sigma}^2) \left(\sum_{i \neq j} a_{ij}b_{ij} + \sum_{i \neq j} a_{ij}b_{ji} \right)$$

Da \mathbf{A}_n und \mathbf{B}_n symmetrisch sind gilt: $\sum_{i \neq j} a_{ij}b_{ji} = \sum_{i \neq j} a_{ij}b_{ij}$ und der obige Ausdruck vereinfacht sich zu

$$E[(\mathbf{Y}'\mathbf{A}\mathbf{Y})(\mathbf{Y}'\mathbf{B}\mathbf{Y})] = \tau_4 \sum_{i=1}^n a_{ii}b_{ii} + [Sp(\boldsymbol{\Sigma})]^2 \sum_{i \neq j} a_{ii}b_{jj} + 2Sp(\boldsymbol{\Sigma}^2) \left(\sum_{i \neq j} a_{ij}b_{ij} \right).$$

Da alle a_{ij} Einträge der Matrix \mathbf{A}_n , sowie alle b_{ij} Einträge der Matrix \mathbf{B}_n sind, wird nur eine Matrixdarstellung der Summanden von \mathbf{A}_n und \mathbf{B}_n benötigt.

Zuerst wird die Summe der Diagonaleinträge von \mathbf{A}_n und \mathbf{B}_n bestimmt:

$$\mathbf{a}_n = \text{diag}\{\mathbf{A}_n\} = (a_{11}, \dots, a_{nn})'$$

$$\mathbf{b}_n = \text{diag}\{\mathbf{B}_n\} = (b_{11}, \dots, b_{nn})' \text{ und es gilt}$$

$$\mathbf{a}'_n \mathbf{b}_n = (a_{11}, \dots, a_{nn}) \cdot (b_{11}, \dots, b_{nn})'$$

$$= \sum_{i=1}^n a_{ii}b_{ii}.$$

Für den zweiten Teil der Summe wird die Summe aller Diagonaleinträge von \mathbf{A}_n und \mathbf{B}_n abzüglich des ersten Teils der Summe benötigt. Dieser lässt sich durch eine Null-Erweiterung mit $\sum_{i=1}^n a_{ii}b_{ii}$ bestimmen:

$$\sum_{i \neq j} a_{ii}b_{jj} = \left(\sum_{i=1}^n a_{ii}b_{ii} + \sum_{i \neq j} a_{ii}b_{jj} \right) - \sum_{i=1}^n a_{ii}b_{ii}$$

$$= Sp(\mathbf{A}_n) Sp(\mathbf{B}_n) - \mathbf{a}'_n \mathbf{b}_n.$$

Für den dritten Teil der Summe wird die Summe der Diagonaleinträge von $\mathbf{A}_n \cdot \mathbf{B}_n$

abzüglich des ersten Teils der Summe benötigt. Das Vorgehen ist dabei analog zu oben:

$$\begin{aligned} 2 \left(\sum_{i \neq j} a_{ij} b_{ji} \right) &= 2 \left(\sum_{i=1}^n \sum_{j=1}^n a_{ij} b_{ji} - \sum_{i=1}^n a_{ii} b_{ii} \right) \\ &= 2Sp(\mathbf{A}_n \mathbf{B}_n) - 2\mathbf{a}'_n \mathbf{b}_n. \end{aligned}$$

Nun ist es möglich den $E[(\mathbf{Y}'\mathbf{A}\mathbf{Y})(\mathbf{Y}'\mathbf{B}\mathbf{Y})]$ in Matrixschreibweise darzustellen:

$$\begin{aligned} E[(\mathbf{Y}'\mathbf{A}\mathbf{Y})(\mathbf{Y}'\mathbf{B}\mathbf{Y})] &= \left(\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 - 2Sp(\boldsymbol{\Sigma}^2) \right) \mathbf{a}'_n \mathbf{b}_n \\ &\quad + [Sp(\boldsymbol{\Sigma})]^2 [Sp(\mathbf{A}_n) Sp(\mathbf{B}_n)] + 2Sp(\boldsymbol{\Sigma}^2) Sp(\mathbf{A}_n \mathbf{B}_n). \end{aligned}$$

Daraus ergibt sich folgende Darstellung für die Kovarianz zweier quadratischer Formen:

$$\begin{aligned} Cov(\mathbf{Y}'\mathbf{A}\mathbf{Y}, \mathbf{Y}'\mathbf{B}\mathbf{Y}) &= E[(\mathbf{Y}'\mathbf{A}\mathbf{Y})(\mathbf{Y}'\mathbf{B}\mathbf{Y})] - E(\mathbf{Y}'\mathbf{A}\mathbf{Y}) E(\mathbf{Y}'\mathbf{B}\mathbf{Y}) \\ &= \left(\tau_4 - [Sp(\boldsymbol{\Sigma})]^2 - 2Sp(\boldsymbol{\Sigma}^2) \right) \mathbf{a}'_n \mathbf{b}_n + 2Sp(\boldsymbol{\Sigma}^2) Sp(\mathbf{A}_n \mathbf{B}_n). \end{aligned}$$

□

Korollar A.2.2 (Darstellung einer Bilinearform (Seite 38))

Für die Zufallsvektoren gelte $\mathbf{X}_k = (X_{k1}, \dots, X_{kd})' = \mathbf{\Gamma}\mathbf{Z}_k + E_k \cdot \mathbf{1}_d + \boldsymbol{\mu}$, $k = 1, \dots, n$, wie in (2.1) mit $\text{Cov}(\mathbf{\Gamma}\mathbf{Z}_k) = \mathbf{\Gamma}\mathbf{\Gamma}' = \mathbf{S}$ und sei $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k + \mathbf{T}\boldsymbol{\mu}$ definiert wie in Definition 2.2 mit $\text{Cov}(\mathbf{Y}_k) = \mathbf{T}\mathbf{S}\mathbf{T}' = \boldsymbol{\Sigma}$. Dann ergibt sich für $k \neq l$ unter $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ das spezielle Resultat von Satz 6.1.4:

$$\begin{aligned} A_{kl} &= \mathbf{X}'_k \mathbf{T}\mathbf{X}_l = \mathbf{Y}'_k \mathbf{Y}_l \\ &= \sum_{i=1}^d \lambda_i Z_{ki} Z_{li}, \end{aligned}$$

wobei die λ_i die Eigenwerte von $\mathbf{T}\mathbf{S}$ sind und die Zufallsvariablen Z_{ki} unabhängig $\forall k = 1, \dots, n, i = 1, \dots, d$.

Beweis:

Unter der Hypothese $H_0 : \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ folgt aus Proposition 2.3.1:

$\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k$ sowie $E_{H_0}(\mathbf{Y}_k) = \mathbf{0}$.

Somit erfüllen die Zufallsvektoren \mathbf{Y}_k , $k = 1, \dots, n$ die Voraussetzungen von Satz 6.1.4 und es kann daher Analog zum Beweis von Satz 6.1.4 vorgegangen werden.

1. Schritt

Die Aussage wird zunächst für $r = d$ bewiesen, d.h. $\det(\mathbf{S}) = |\mathbf{S}| \neq 0$.

Da $\mathbf{\Gamma}$ frei wählbar ist solange die Bedingung $\mathbf{\Gamma}\mathbf{\Gamma}' = \mathbf{S}$ erfüllt wird, kann für diesen Fall $\mathbf{\Gamma} = \mathbf{S}^{1/2}\mathbf{P}$ gesetzt werden, wobei $\mathbf{S}^{1/2}$ die Wurzel der Kovarianzmatrix des Zufallsvektors \mathbf{X} ist und \mathbf{P} die orthogonale Matrix, welche zum diagonalisieren der Matrix $\mathbf{S}^{1/2}\mathbf{T}\mathbf{S}^{1/2}$ benötigt wird.

Dann gilt:

$$\begin{aligned} \mathbf{T}\mathbf{X}_k &= \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k = \mathbf{T}\mathbf{S}^{1/2}\mathbf{P}\mathbf{Z}_k \text{ und} \\ \mathbf{T}\mathbf{X}_l &= \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_l = \mathbf{T}\mathbf{S}^{1/2}\mathbf{P}\mathbf{Z}_l \end{aligned}$$

mit

$$\mathbf{\Gamma}\mathbf{\Gamma}' = \mathbf{S}^{1/2}\mathbf{P}\mathbf{P}'\mathbf{S}^{1/2} = \mathbf{S}^{1/2}\mathbf{S}^{1/2} = \mathbf{S}.$$

Der weitere Beweis ist dann analog zum Darstellungssatz 6.1.4 und wird im Folgenden kurz aufgezeigt.

Es gilt dann: $|\mathbf{S}| \neq 0$ und $\mathbf{S} = \mathbf{S}' \Rightarrow \exists \mathbf{S}^{1/2}$ und $\mathbf{S}^{-1/2}$, symmetrisch und invertierbar $\Rightarrow \mathbf{S}^{1/2} \mathbf{S}^{-1/2} = \mathbf{I}_d$.

$$\begin{aligned} A_{kl} &= \mathbf{X}'_k \mathbf{T} \mathbf{X}_l = \mathbf{Y}'_k \mathbf{Y}_l = \left(\mathbf{S}^{\frac{1}{2}} \mathbf{P} \mathbf{Z}_k \right)' \mathbf{T} \left(\mathbf{S}^{\frac{1}{2}} \mathbf{P} \mathbf{Z}_l \right) \\ &= \left(\mathbf{S}^{\frac{1}{2}} \mathbf{P} \mathbf{Z}_k \right)' \mathbf{S}^{-1/2} \mathbf{S}^{1/2} \mathbf{T} \mathbf{S}^{1/2} \mathbf{S}^{-1/2} \left(\mathbf{S}^{\frac{1}{2}} \mathbf{P} \mathbf{Z}_l \right) \\ &= \left(\mathbf{S}^{\frac{1}{2}} \mathbf{P} \mathbf{Z}_k \right)' \mathbf{S}^{-1/2} \mathbf{P}' \mathbf{P} \mathbf{S}^{1/2} \mathbf{T} \mathbf{S}^{1/2} \mathbf{P}' \mathbf{P} \mathbf{S}^{-1/2} \left(\mathbf{S}^{\frac{1}{2}} \mathbf{P} \mathbf{Z}_l \right) \end{aligned}$$

Da $\mathbf{S}^{1/2}$ und \mathbf{T} symmetrisch sind, gilt auch, dass $\mathbf{S}^{1/2} \mathbf{T} \mathbf{S}^{1/2}$ symmetrisch ist. Damit folgt aus dem Satz über die Hauptachsentransformation: Es existiert eine orthogonale Matrix \mathbf{P} , sodass sich $\mathbf{P} \mathbf{S}^{1/2} \mathbf{T} \mathbf{S}^{1/2} \mathbf{P}' = \text{diag} \{ \lambda'_1, \dots, \lambda'_d \} = \Delta$ darstellen lässt. Die λ'_i sind hierbei die Eigenwerte von $\mathbf{S}^{1/2} \mathbf{T} \mathbf{S}^{1/2}$.

$$\begin{aligned} \mathbf{X}'_k \mathbf{T} \mathbf{X}_l &= \left(\mathbf{P} \mathbf{S}^{-1/2} \mathbf{P}' \mathbf{Z}_k \right)' \underbrace{\mathbf{P} \mathbf{S}^{1/2} \mathbf{T} \mathbf{S}^{1/2} \mathbf{P}'}_{\Delta} \left(\mathbf{P} \mathbf{S}^{-1/2} \mathbf{P}' \mathbf{Z}_l \right) \\ &= \left(\mathbf{P} \mathbf{S}^{-1/2} \mathbf{S}^{1/2} \mathbf{P}' \mathbf{Z}_k \right)' \Delta \left(\mathbf{P} \mathbf{S}^{-1/2} \mathbf{S}^{1/2} \mathbf{P}' \mathbf{Z}_l \right) \\ &= \mathbf{Z}'_k \Delta \mathbf{Z}_l = \sum_{i=1}^d \lambda_i Z_{ki} Z_{li} \end{aligned}$$

Im Beweis von Lemma 6.1.4 wurde allgemein gezeigt, dass $\lambda'_i = \lambda_i$ gilt, mit λ_i sind die Eigenwerte von $\mathbf{T} \mathbf{S}$ (siehe dafür Werner (2004) oder Brunner (2010)).

2. Fall

Sei nun \mathbf{S} singulär mit $r(\mathbf{S}) = r < d$:

Für diesen Fall wird $\mathbf{\Gamma} = \mathbf{P}' \left(\begin{array}{c|c} \mathbf{S}_r^{1/2} \mathbf{P}'_{r^*} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right)$ gewählt, mit \mathbf{P} und \mathbf{P}_{r^*} als orthogonale Matrizen mit $\mathbf{P} \mathbf{S} \mathbf{P}' = \left(\begin{array}{c|c} \mathbf{S}_r & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right)$ und \mathbf{S}_r sei die Diagonalmatrix, welche die von Null verschiedenen Eigenwerte der Kovarianzmatrix \mathbf{S} enthält. Zur besseren Übersicht wird $\tilde{\mathbf{\Sigma}}^{1/2} = \left(\begin{array}{c|c} \mathbf{S}_r^{1/2} \mathbf{P}'_{r^*} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right)$ definiert, sowie anstelle der Zufallsvektoren \mathbf{X}_k der Fehlerterm $\mathbf{\Gamma} \mathbf{Z}_k$ betrachtet. Dann gilt:

$$\begin{aligned}
 \Gamma Z_k &= P' \left(\begin{array}{c|c} \mathbf{S}_r^{1/2} \mathbf{P}'_{r^*} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right) Z_k \\
 &= P' \tilde{\Sigma}^{1/2} Z_k \text{ und} \\
 \Gamma Z_l &= P' \left(\begin{array}{c|c} \mathbf{S}_r^{1/2} \mathbf{P}'_{r^*} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right) Z_l \\
 &= P' \tilde{\Sigma}^{1/2} Z_l,
 \end{aligned}$$

mit

$$\begin{aligned}
 \Gamma \Gamma' &= P' \left(\begin{array}{c|c} \mathbf{S}_r^{1/2} \mathbf{P}'_{r^*} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right) \left(\begin{array}{c|c} \mathbf{P}_{r^*} \mathbf{S}_r^{1/2} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right) P \\
 &= P' \left(\begin{array}{c|c} \mathbf{S}_r^{1/2} \mathbf{P}'_{r^*} \mathbf{P}_{r^*} \mathbf{S}_r^{1/2} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right) P = P' \left(\begin{array}{c|c} \mathbf{S}_r & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right) P = \mathbf{S}.
 \end{aligned}$$

Da \mathbf{S} symmetrisch ist, folgt aus der Hauptachsentransformation, dass eine orthogonale Matrix \mathbf{P} existiert, so dass $\mathbf{S}_r = \text{diag}\{\nu_1, \dots, \nu_r\}$ ist, wobei die ν_i , $i = 1, \dots, r$, die von 0 verschiedenen Eigenwerte von \mathbf{S} sind. Damit ist $|\mathbf{S}_r| \neq 0$.

Setzen nun

$$\begin{aligned}
 \mathbf{V} &= P \Gamma Z_k = P P' \tilde{\Sigma}^{1/2} Z_k = \tilde{\Sigma}^{1/2} (Z'_{k1} | Z'_{k2})' = (\mathbf{V}'_1 | \mathbf{V}'_2)' \text{ und} \\
 \mathbf{W} &= P \Gamma Z_l = P P' \tilde{\Sigma}^{1/2} Z_l = \tilde{\Sigma}^{1/2} (Z'_{l1} | Z'_{l2})' = (\mathbf{W}'_1 | \mathbf{W}'_2)
 \end{aligned}$$

mit

$$\begin{aligned}
 \mathbf{V}_1 &= (V_1, \dots, V_r)', \\
 \mathbf{V}_2 &= (V_{r+1}, \dots, V_d)' \text{ und} \\
 \mathbf{W}_1 &= (W_1, \dots, W_r)', \\
 \mathbf{W}_2 &= (W_{r+1}, \dots, W_d)',
 \end{aligned}$$

sowie entsprechend

$$\begin{aligned}
 \mathbf{Z}_{k1} &= (Z_{k1}, \dots, Z_{kr})', \\
 \mathbf{Z}_{k2} &= (Z_{k,r+1}, \dots, Z_{kd})' \text{ und}
 \end{aligned}$$

$$\begin{aligned}\mathbf{Z}_{l1} &= (Z_{l1}, \dots, Z_{lr})', \\ \mathbf{Z}_{l2} &= (Z_{l,r+1}, \dots, Z_{ld})' .\end{aligned}$$

Dann folgt weiter

$$\text{Cov}(\mathbf{P}\Gamma\mathbf{Z}_k) = \text{Cov}(\mathbf{P}\Gamma\mathbf{Z}_k) = \left(\tilde{\Sigma}^{1/2}\right) \mathbf{I}_d \left(\tilde{\Sigma}^{1/2}\right)' = \left(\begin{array}{c|c} \mathbf{S}_r & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} \end{array}\right).$$

Damit ist $\text{Cov}(\mathbf{V}_2) = \text{Cov}(\mathbf{W}_2) = \mathbf{0}$. Da $E(\Gamma\mathbf{Z}_k) = E(\Gamma\mathbf{Z}_l) = \mathbf{0}$ gilt, folgt daraus $E(\mathbf{P}\Gamma\mathbf{Z}_k) = E(\mathbf{P}\Gamma\mathbf{Z}_l) = \mathbf{0}$ und damit $E(\mathbf{V}_2) = E(\mathbf{W}_2) = \mathbf{0}$. Daraus ergibt sich, dass $\mathbf{V}_2 = \mathbf{0}$ f.s. und $\mathbf{W}_2 = \mathbf{0}$ f.s.

Weiter folgt dann für die Bilinearform mit $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\Gamma\mathbf{X}_k$:

$$\begin{aligned}A_{kl} &= \mathbf{X}'_k \mathbf{T} \mathbf{X}_l = \mathbf{Y}'_k \mathbf{Y}_l = \left(\mathbf{P}' \tilde{\Sigma}^{1/2} (\mathbf{Z}'_{k1} | \mathbf{Z}'_{k2})'\right)' \mathbf{T} \left(\mathbf{P}' \tilde{\Sigma}^{1/2} (\mathbf{Z}'_{l1} | \mathbf{Z}'_{l2})'\right) \\ &= \left(\mathbf{P}' \tilde{\Sigma}^{1/2} (\mathbf{Z}'_{k1} | \mathbf{Z}'_{k2})'\right)' \underbrace{\mathbf{P}' \mathbf{P} \mathbf{T} \mathbf{P}' \mathbf{P}}_{\tilde{\mathbf{T}}} \left(\mathbf{P}' \tilde{\Sigma}^{1/2} (\mathbf{Z}'_{l1} | \mathbf{Z}'_{l2})'\right) \\ &= \left(\tilde{\Sigma}^{1/2} (\mathbf{Z}'_{k1} | \mathbf{Z}'_{k2})'\right)' \tilde{\mathbf{T}} \left(\tilde{\Sigma}^{1/2} (\mathbf{Z}'_{l1} | \mathbf{Z}'_{l2})'\right) \\ &= \left((\mathbf{V}'_1 | \mathbf{V}'_2)'\right)' \left(\begin{array}{c|c} \tilde{\mathbf{T}}_{11} & \tilde{\mathbf{T}}_{12} \\ \hline \tilde{\mathbf{T}}_{21} & \tilde{\mathbf{T}}_{22} \end{array}\right) \left((\mathbf{W}'_1 | \mathbf{W}'_2)'\right) = \mathbf{V}'_1 \tilde{\mathbf{T}}_{11} \mathbf{W}_1 \quad f.s. \\ &= \left(\mathbf{S}_r^{1/2} \mathbf{P}'_{r*} \mathbf{Z}_{k1}\right)' \tilde{\mathbf{T}}_{11} \left(\mathbf{S}_r^{1/2} \mathbf{P}'_{r*} \mathbf{Z}_{l1}\right) \quad f.s..\end{aligned}$$

Abschließend kann auf die Bilinearform $\mathbf{V}'_1 \tilde{\mathbf{T}}_{11} \mathbf{W}_1$ der erste Teil dieses Darstellungssatzes (Satz A.2.2) angewendet werden, da die zugehörige Kovarianzmatrix \mathbf{S}_r nicht singulär ist.

$$\begin{aligned}\mathbf{V}'_1 \tilde{\mathbf{T}}_{11} \mathbf{W}_1 &= \left(\mathbf{S}_r^{1/2} \mathbf{P}'_{r*} \mathbf{Z}_{k1}\right)' \tilde{\mathbf{T}}_{11} \left(\mathbf{S}_r^{1/2} \mathbf{P}'_{r*} \mathbf{Z}_{l1}\right) \\ &= \mathbf{Z}'_{k1} \underbrace{\mathbf{P}_{r*} \mathbf{S}_r^{1/2} \tilde{\mathbf{T}}_{11} \mathbf{S}_r^{1/2} \mathbf{P}'_{r*}}_{\Delta} \mathbf{Z}_{l1} \\ &= \mathbf{Z}'_{k1} \Delta \mathbf{Z}_{l1} = \sum_{i=1}^r \tilde{\lambda}_i Z_{ki} Z_{li}\end{aligned}$$

Die $\tilde{\lambda}_i$ sind, nach Lemma 6.1.4, die Eigenwerte von $\mathbf{S}_r^{1/2} \tilde{\mathbf{T}}_{11} \mathbf{S}_r^{1/2}$ und $\tilde{\mathbf{T}}_{11} \mathbf{S}_r$ sowie die von Null verschiedenen Eigenwerte von $\mathbf{T}\mathbf{S}$.

□

Satz A.2.3 (Schätzer Eigenschaften (Seite 43))

Seien B_1 und B_2 definiert wie in Definition 6.1.2 und erfüllen die Zufallsvektoren \mathbf{X}_k , $k = 1, \dots, n$, die Modellannahmen aus (2.1). Weiterhin sei $\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k + \mathbf{T}\boldsymbol{\mu}$ definiert wie in Proposition 2.3.1 mit $\text{Cov}(\mathbf{Y}_k) = \mathbf{T}\mathbf{S}\mathbf{T} = \boldsymbol{\Sigma}$.

Somit gilt unter $H_0: \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$: B_1 und B_2

1. sind erwartungstreue Schätzer für $[\text{Sp}(\boldsymbol{\Sigma})]^2$ und $\text{Sp}(\boldsymbol{\Sigma}^2)$.
2. sind konsistent im Sinne von Definition A.3.2.
3. sind dimensionsstabil im Sinne von Definition A.3.4.

Beweis:

Unter $H_0: \mathbf{T}\boldsymbol{\mu} = \mathbf{0}$ folgt aus Proposition 2.3.1:

$\mathbf{Y}_k = \mathbf{T}\mathbf{X}_k = \mathbf{T}\mathbf{\Gamma}\mathbf{Z}_k$ sowie $E_{H_0}(\mathbf{Y}_k) = \mathbf{0}$.

Damit ergeben sich die folgenden Resultate:

1.

Die Erwartungstreue der Schätzer kann mit Hilfe von Lemma 6.1.3 (Momente I) ohne Verwendung der Modellannahmen (2.1) auf Seite 5 nachgewiesen werden:

$$B_1 = \frac{1}{n(n-1)} \cdot \sum_{k \neq l} A_k \cdot A_l \qquad B_2 = \frac{1}{n \cdot (n-1)} \cdot \underbrace{\sum_{k=1}^n \sum_{l=1}^n}_{k \neq l} A_{kl}^2$$

mit

$$E_{H_0}(B_1) = \frac{1}{n(n-1)} \sum_{k \neq l} E(A_k A_l) = \frac{n(n-1)}{n(n-1)} \cdot [\text{Sp}(\boldsymbol{\Sigma})]^2 = [\text{Sp}(\boldsymbol{\Sigma})]^2 \text{ und}$$

$$E_{H_0}(B_2) = \frac{1}{n(n-1)} \sum_{k \neq l} E(A_{kl}^2) = \frac{n(n-1)}{n(n-1)} \cdot \text{Sp}(\boldsymbol{\Sigma}^2) = \text{Sp}(\boldsymbol{\Sigma}^2).$$

2. und 3.

Konsistenz und Dimensionsstabilität werden jeweils für B_1 und B_2 getrennt gezeigt.

Um diese Eigenschaften zu zeigen, wird zuerst eine Abschätzung für die Varianz des Schätzers B_1 benötigt, welche sich mit Hilfe der vorangegangenen Lemma 6.1.3 und 6.1.8 leicht bestimmen lässt.

$$\begin{aligned}
Var(B_1) &= Var \left(\frac{1}{n(n-1)} \sum_{k \neq l} A_k A_l \right) \\
&= \frac{1}{n^2(n-1)^2} \left[E \left(\left[\sum_{k \neq l} A_k A_l \right]^2 \right) - \left[E \left(\sum_{k \neq l} A_k A_l \right) \right]^2 \right] \\
&= \frac{1}{n^2(n-1)^2} \left[E \left(\sum_{k \neq l} \sum_{r \neq s} A_k A_l A_r A_s \right) \right] - [E(A_k A_l)]^2 \\
&= \frac{1}{n^2(n-1)^2} \left[\sum_{k \neq l} \sum_{r \neq s} E(A_k A_l A_r A_s) \right] - [Sp(\Sigma)]^4 \\
&= \frac{1}{n^2(n-1)^2} \left[2 \sum_{k \neq l} E(A_k^2 A_l^2) + 4 \sum_{k \neq l} \sum_{\neq s} E(A_k^2) E(A_l A_s) \right] \\
&\quad + \frac{1}{n^2(n-1)^2} \left[\sum_{k \neq l} \sum_{\neq r \neq s} E(A_k A_l A_r A_s) \right] - [Sp(\Sigma)]^4 \\
&\leq \frac{1}{n(n-1)} \left[\underbrace{2 \left(\gamma Sp(\Sigma^2) + [Sp(\Sigma)]^2 \right)^2}_A \right] \\
&\quad + \frac{1}{n(n-1)} \left[4(n-2) \underbrace{\left(\gamma Sp(\Sigma^2) [Sp(\Sigma)]^2 + [Sp(\Sigma)]^4 \right)}_B \right] \\
&\quad + \frac{(n-2)(n-3)}{n(n-1)} [Sp(\Sigma)]^4 - \frac{n(n-1)}{n(n-1)} [Sp(\Sigma)]^4
\end{aligned}$$

Die Abschätzung von A und B mit Hilfe von Lemma 6.1.3 (5) (Momente I) ergibt:

A

$$\begin{aligned}
\left(\gamma Sp(\Sigma^2) + [Sp(\Sigma)]^2 \right)^2 &\leq \left(\gamma [Sp(\Sigma)]^2 + [Sp(\Sigma)]^2 \right)^2 \\
&= [Sp(\Sigma)]^4 (\gamma + 1)^2,
\end{aligned}$$

B

$$\gamma Sp(\boldsymbol{\Sigma}^2) [Sp(\boldsymbol{\Sigma})]^2 + [Sp(\boldsymbol{\Sigma})]^4 \leq [Sp(\boldsymbol{\Sigma})]^4 \cdot (\gamma + 1).$$

Durch diese Abschätzungen folgt für die Varianz von B_1 :

$$\begin{aligned} Var(B_1) &\leq \frac{1}{n(n-1)} \left[2 \cdot A + 4(n-2) \cdot B - (4n-6) \cdot [Sp(\boldsymbol{\Sigma})]^4 \right] \\ &\leq \frac{[Sp(\boldsymbol{\Sigma})]^4}{n(n-1)} \left[2(\gamma+1)^2 + 4(n-2)(\gamma+1) - (4n-6) \right]. \end{aligned}$$

Für die Dimensionsstabilität wird gezeigt, dass $Var\left(\frac{B_1}{[Sp(\boldsymbol{\Sigma})]^2}\right)$ eine von der Dimension unabhängige Nullfolge als obere Schranke besitzt (siehe auch Definition A.3.4).

$$\begin{aligned} Var\left(\frac{B_1}{[Sp(\boldsymbol{\Sigma})]^2}\right) &= \frac{Var(B_1)}{[Sp(\boldsymbol{\Sigma})]^4} \\ &\leq \frac{[Sp(\boldsymbol{\Sigma})]^4 \cdot \left[2 \cdot (\gamma+1)^2 + 4(n-2) \cdot (\gamma+1) - (4n-6) \right]}{n(n-1) [Sp(\boldsymbol{\Sigma})]^4} \\ &= \frac{\left[2 \cdot (\gamma+1)^2 + 4(n-2) \cdot (\gamma+1) - (4n-6) \right]}{n(n-1)} \\ &\leq \frac{2}{n(n-1)} \cdot (\gamma+1)^2 + 4 \frac{1}{n} \cdot (\gamma+1) - \frac{2 \cdot (2 - \frac{3}{n})}{(n-1)} = C_1(n) \end{aligned}$$

$C_1(n)$ ist unabhängig von der Wahl der Dimension d und für $n \rightarrow \infty$ geht $C_1(n) \rightarrow 0$. Damit ist die Dimensionsstabilität im Sinne von Definition A.3.4 erfüllt. Da $C_1(n)$ eine Nullfolge ist, folgt auch direkt die Konsistenz im Sinne von Definition A.3.2.

Analog wird auch für den Schätzer B_2 zuerst eine Abschätzung für die Varianz mit Hilfe des Lemma 6.1.8 bestimmt.

$$\begin{aligned} Var(B_2) &= Var\left(\frac{1}{n(n-1)} \sum_{k \neq l} A_{kl}^2\right) \\ &= \frac{1}{n^2(n-1)^2} \left[E\left(\sum_{k \neq l} \sum_{lr \neq s} A_{kl}^2 A_{rs}^2\right) - \left[E\left(\sum_{k \neq l} A_{kl}^2\right) \right]^2 \right] \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{n^2(n-1)^2} \left[\sum_{k \neq l} \sum_{lr \neq s} E(A_{kl}^2 A_{rs}^2) \right] - [Sp(\Sigma^2)]^2 \\
&\leq \frac{(4n-6)}{n(n-1)} E(A_{kl}^4) + \frac{n(n-1)(n-2)(n-3)}{n^2(n-1)^2} [Sp(\Sigma^2)]^2 - [Sp(\Sigma^2)]^2 \\
&= \frac{(4n-6)}{n(n-1)} \left[E(A_{kl}^4) - [Sp(\Sigma^2)]^2 \right] \\
&\leq \frac{(4n-6)}{n(n-1)} \left(\gamma Sp(\Sigma^4) + 3 [Sp(\Sigma^2)]^2 - [Sp(\Sigma^2)]^2 \right) \\
&= \frac{(4n-6)}{n(n-1)} \left(\gamma Sp(\Sigma^4) + 2 [Sp(\Sigma^2)]^2 \right)
\end{aligned}$$

Aus der Cauchy-Schwarz Ungleichung und der Darstellung von $Sp(\Sigma)$ als Summe der Eigenwerte λ_i , $i = 1, \dots, d$, von Σ folgt:

$$Sp(\Sigma^4) = \sum_{i=1}^d \lambda_i^4 = \sum_{i=1}^d \tilde{\lambda}_i^2 \leq \left(\sum_{i=1}^d \tilde{\lambda}_i \right)^2 = \left(\sum_{i=1}^d \lambda_i^2 \right)^2 = [Sp(\Sigma^2)]^2$$

mit $\tilde{\lambda}_i = \lambda_i^2$, $i = 1, \dots, d$.

Nun lässt sich auch die Dimensionsstabilität von B_2 nachweisen:

$$\begin{aligned}
Var \left(\frac{B_2}{Sp(\Sigma^2)} \right) &= \frac{Var(B_2)}{[Sp(\Sigma^2)]^2} \leq \frac{(4n-6)}{n(n-1)} \cdot \frac{\left(\gamma Sp(\Sigma^4) + 2 [Sp(\Sigma^2)]^2 \right)}{[Sp(\Sigma^2)]^2} \\
&\leq \frac{(4n-6)}{n(n-1)} \cdot \frac{\left(\gamma [Sp(\Sigma^2)]^2 + 2 [Sp(\Sigma^2)]^2 \right)}{[Sp(\Sigma^2)]^2} \\
&= \frac{(4n-6)}{n(n-1)} \cdot [\gamma + 2] = C_2(n).
\end{aligned}$$

$C_2(n)$ ist unabhängig von der Wahl der Dimension d und für $n \rightarrow \infty$ geht $C_2(n) \rightarrow 0$. Damit ist die Dimensionsstabilität im Sinne von Definition A.3.4 erfüllt. Da $C_2(n)$ eine Nullfolge ist, folgt auch die Konsistenz im Sinne von Definition A.3.2.

□

A.3. Verwendete Definitionen, Lemmata und Sätze

Definition A.3.1 (Konsistenz 1)

Eine Folge von Schätzern $\hat{\theta}_n$ heißt konsistent für θ , falls $\hat{\theta}_n - \theta \xrightarrow{P} 0$ für festes θ (d.h. $\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| > \epsilon) = 0 \forall \epsilon > 0$).

Die Konsistenz ist eine Minimaleigenschaft, die Schätzer mindestens erfüllen sollten. Sie impliziert, dass sich der Schätzer mit wachsendem Stichprobenumfang n dem zu schätzenden Parameter immer besser annähert.

Definition A.3.2 (Konsistenz 2)

Ein Feld von Schätzern $\hat{\theta}_{n,d}$ heißt konsistent für θ_d , falls $\hat{\theta}_{n,d} - \theta_d \xrightarrow{P} 0$ für festes d , somit auch für festes θ_d (d.h. $\lim_{n \rightarrow \infty} P(|\hat{\theta}_{n,d} - \theta_d| > \epsilon) = 0 \forall \epsilon > 0$).

Lemma A.3.3 Ein Feld von asymptotisch erwartungstreuen Schätzern $\hat{\theta}_{n,d}$ ist konsistent für $\theta_d > 0$, falls gilt:

$$\lim_{n \rightarrow \infty} \left\{ \frac{\text{Var}(\hat{\theta}_{n,d})}{\theta_d^2} \right\} = 0 \quad \forall d < \infty, d \text{ fest.}$$

Beweis: Siehe Werner (2004), S. 11-12

□

Definition A.3.4 (Dimensionsstabilität)

Ein Feld von Schätzern $\hat{\theta}_{n,d}$ heißt dimensionsstabil, wenn $\forall d \geq 1, n \geq 1$ gilt:

1. $|E\left(\frac{\hat{\theta}_{n,d}}{\theta_d} - 1\right)| \leq A(n) < \infty$, wobei A nicht von der Dimension der Messwiederholungen d abhängt.
2. $\text{Var}\left(\frac{\hat{\theta}_{n,d}}{\theta_d}\right) = \frac{1}{\theta_d^2} \text{Var}(\hat{\theta}_{n,d}) \leq B(n) < \infty$, wobei B nicht von der Dimension d abhängt.

Die Definitionen bezüglich eines Feldes von Schätzern stammen aus Werner (2004).

Definition A.3.5 (Asymptotische Äquivalenz)

Zwei Folgen von Zufallsvariablen X_n und Y_n sind asymptotisch äquivalent (in Zeichen $X_n \doteq Y_n$), wenn $X_n - Y_n \xrightarrow{P} 0$ für $n \rightarrow \infty$ gilt.

Unter Verwendung der tschebyschewschen Ungleichung (siehe z.B. Krenzel, 2002, Satz 3.15, S. 57) ist es zumeist einfacher, das stärkere Resultat $E[(X_n - Y_n)^2] \rightarrow 0$ zu zeigen. Asymptotische Äquivalenz impliziert asymptotische Verteilungsgleichheit.

Satz A.3.6 (Lancaster)

Sei $\mathbf{X} = (X_1, \dots, X_n)'$ ein Zufallsvektor mit $E(\mathbf{X}) = \boldsymbol{\mu} = (\mu_1, \dots, \mu_n)'$ und $\mathbf{V} = \text{Var}(\mathbf{X})$. Ferner sei $\mathbf{A} = \mathbf{A}'$ und $\text{Sp}(\cdot)$ bezeichne die Spur einer Matrix. Dann gilt:

$$E(\mathbf{X}'\mathbf{A}\mathbf{X}) = \text{Sp}(\mathbf{A}\mathbf{V}) + \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}.$$

Beweis: Siehe Mathai und Provost (1992), S. 53

□

Satz A.3.7 (Varianz einer quadratischen Form (Atiqullah))

Die Zufallsvariablen X_i , $i = 1, \dots, n$, seien unabhängig und $\mathbf{X} = (X_1, \dots, X_n)'$ bezeichne den Vektor der Zufallsvariablen und $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)'$ den Vektor der Erwartungswerte. Ferner sei $\text{Cov}(\mathbf{X}) = \sigma^2 \mathbf{I}_n$.

Weiterhin wird angenommen, dass die X_i identische dritte und vierte Momente besitzen, d.h. $E[(X_i - \mu_i)^3] = \mu_3$ und $E[(X_i - \mu_i)^4] = \mu_4$, $i = 1, \dots, n$. Sei $\mathbf{A}' = \mathbf{A}_{n \times n}$ und bezeichne $\mathbf{a} = \text{diag}\{\mathbf{A}\}$ den Vektor der Diagonalelemente von \mathbf{A} . Dann ist

$$\text{Var}(\mathbf{X}'\mathbf{A}\mathbf{X}) = (\mu_4 - 3\sigma^4) \mathbf{a}'\mathbf{a} + 2\sigma^4 \text{Sp}(\mathbf{A}^2) + 4\sigma^4 \boldsymbol{\mu}'\mathbf{A}^2\boldsymbol{\mu} + 4\mu_3 \boldsymbol{\mu}'\mathbf{A}\mathbf{a}.$$

Beweis: Siehe Atiqullah (1962) und Brunner (2010)

□

Lemma A.3.8 Sei $X_k \sim N(\boldsymbol{\mu}, \mathbf{V})$, $\mathbf{V} \in \mathbb{R}^{n \times n}$. Seien $\mathbf{L}, \mathbf{M} \in \mathbb{R}^{n \times k}$ Matrizen mit Spaltenrang $k \leq n$. Dann sind $\mathbf{L}'\mathbf{X}$ und $\mathbf{M}'\mathbf{X}$ unabhängig, falls $\text{Cov}(\mathbf{L}'\mathbf{X}, \mathbf{M}'\mathbf{X}) = \mathbf{0}$ ist.

Beweis: Siehe Brunner (2010)

□

Satz A.3.9 (Craig und Sakamoto)
Sei $\mathbf{X} \sim N(\boldsymbol{\mu}, \mathbf{V})$, $\mathbf{A} = \mathbf{A}'$ p.s.d., $\mathbf{B} = \mathbf{B}'$ p.s.d. und \mathbf{b} ein konstanter Vektor.
Dann gilt:

1. $\mathbf{X}'\mathbf{A}\mathbf{X}$ und $\mathbf{X}'\mathbf{B}\mathbf{X}$ bzw. $\mathbf{X}'\mathbf{A}\mathbf{X}$ und $\mathbf{B}\mathbf{X}$ sind stochastisch unabhängig, falls $\mathbf{B}\mathbf{V}\mathbf{A} = \mathbf{0}$ ist,
2. $\mathbf{A}\mathbf{X}$ und $\mathbf{b}'\mathbf{X}$ sind stochastisch unabhängig, falls $\mathbf{b}'\mathbf{V}\mathbf{A} = \mathbf{0}$ ist.

Beweis: Siehe Craig (1943) oder Brunner (2010)

□

Literaturverzeichnis

- [1] Ahmad, R. M. (2008). Analysis of High Dimensional Repeated Measures Designs: The One- and Two-Sample Test Statistics. Dissertation, Universität Göttingen.
- [2] Ahmad, R. M., Werner, C. und Brunner, E. (2008). Analysis of high-dimensional repeated measures designs: The one sample case. *Computational Statistics and Data Analysis* **53**, 416-427.
- [3] Akritas, M. G., Wang, H. (2010). Inference from Heteroscedastic Functional Data, Part I: Identical Sampling Points. *Journal of Nonparametric Statistics* **22**, 149-168.
- [4] Atiqullah, M. (1962). The estimation of residual variance in quadratically balanced least-squares problems and the robustness of the F-test. *Biometrika* **49**, 1 and 2, 83.
- [5] Bai, Z. und Saranadasa, H. (1996). Effect of highdimension: by an example of a two sample problem. *Statistica Sinica* **6**, 311-329.
- [6] Becker, B. (2010). Test für hochdimensionale Messwiederholungen mit unbekanntem Kovarianzmatrizen. Diplomarbeit, Universität Göttingen.
- [7] Boysen, L. (2002). Analyse von intra-individuellen Effekten bei longitudinalen Daten. Diplomarbeit, Institut für Mathematische Stochastik, Göttingen.
- [8] Brunner, E. (2009). Repeated Measures under Non-Sphericity. Proceedings of the 6-th St. Petersburg Workshop on Simulation. P. 605-610.
- [9] Brunner, E. (2010). Angewandte Statistik II. Vorlesungsskript, Universität Göttingen.
- [10] Brunner, E., Domhof, S. und Langer, F. (2002). *Nonparametric Analysis of Longitudinal Data in Factorial Experiments*. Wiley, New York.
- [11] Brunner, E. und Munzel, M. (2002). *Nichtparametrische Datenanalyse: Unverbundene Stichproben*. Springer, Berlin.
- [12] Box, G. E. P. (1954). Some Theorems on Quadratic Forms Applied in the Study of Analysis of Variance Problems, I. Effect of Inequality of Variance in the One-Way Classification. *The Annals of Mathematical Statistics* **25**, 290-302.

- [13] Box, G. E. P. (1954). Some Theorems on Quadratic Forms Applied in the Study of Analysis of Variance Problems, II. Effects of Inequality of Variance and of Correlation Between Errors in the Two-Way Classification. *The Annals of Mathematical Statistics* **25**, 484-498.
- [14] Craig, A. T. (1943). Note on the Independence of Certain Quadratic Forms. *The Annals of Mathematical Statistics* **14**, No. 2 , 195-197.
- [15] Dehling, H., Haupt, B. (2004). *Einführung in die Wahrscheinlichkeitstheorie und Statistik*. Springer, Berlin.
- [16] Ferguson, T. S. (1996). *A Course in Large Sample Theory*. Chapman, New York.
- [17] Geisser, S. and Greenhouse, S. W. (1958). An Extension of Box's Results on the Use of the F Distribution in Multivariate Analysis. *The Annals of Mathematical Statistics* **29**, No. 3, 885-891.
- [18] Hotelling H. (1931). The generalization of Student's ratio. *The Annals of Mathematical Statistics* **2**, 360-378.
- [19] Johnson, N. L., Kotz, S. und Balakrishnan, N. (1994). *Continuous Univariate Distributions Volume 1*. Wiley, New York.
- [20] Johnson, N. L., Kotz, S. und Balakrishnan, N. (2000). *Distributions in Statistics: Continuous Multivariate Distributions*. Wiley, New York.
- [21] Krenzel, U. (2002). *Einführung in die Wahrscheinlichkeitstheorie und Statistik*. Vieweg, Braunschweig, Wiesbaden, sechste Auflage.
- [22] Mattai, A.M. and Provost, S. B. (1992). *Quadratic forms in random variables: Theory and applications*. Marcel Dekker, INC., New York
- [23] Ohtaki, M. (1990). Some estimators of covariance matrix in multivariate nonparametric regression and their applications. *Hiroshima Math. J.* **20**, 63-91.
- [24] Srivastava, M. S. (2009). A Review of Multivariate Theory For High Dimensional Data with Fewer Observations. *Advances in Multivariate Statistical Methods* **9**, Editor Ashis SenGupta, 25-52
- [25] Strange, K. (1970). *Angewandte Statistik I*. Springer, Berlin
- [26] Werner, C. (2004). Dimensionsstabile Approximation für Verteilungen von quadratischen Formen im Repeated-Measures-Design. Diplomarbeit, Universität Göttingen.