

Theoretische Grundlagen der partiellen kleinsten Quadrate

Dissertation

zur Erlangung des akademischen Grades

doctor rerum naturalium

(Dr. rer. nat.)

von B.Sc./D.H.S. Hayan Hasan

geb. am 12.04.1976 Safita-Syrien

genehmigt durch die Fakultät für Mathematik
der Otto-von-Guericke-Universität Magdeburg

Gutachter: Prof.Dr.Rainer Schwabe

Prof.Dr.Rainer Schlittgen

eingereicht am: 07.01.2014

Verteidigung am: 06.06.2014

Inhaltsverzeichnis

1. Allgemeine Grundlagen der mathematischen Statistik	15
1.1. Statistisches Experiment	15
1.2. Schätzproblem	17
1.2.1. Optimalitätskriterien	18
1.2.2. Eigenschaften des UMVUE-Schätzers	19
1.3. Verteilungsmodelle	23
1.4. Lineares Normalverteilungsmodell	24
1.5. Modell mit multivariaten, u.i.v. normalverteilten ZV'en	25
1.6. Autoregressives Modell	28
2. Assoziationsstrukturen	29
2.1. Theorie der Regressionsmodelle	29
2.2. Kanonische Korrelation	32
2.3. Normalverteilungsmodell	35
3. Einige Aspekte der Theorie der Vorhersage	37
3.1. Allgemeine Resultate	37
3.2. Optimalitätskonzepte	40
3.3. Die beste erwartungstreue Vorhersage	42
3.4. Eigenschaften des mittleren quadratischen Fehlers	49
3.5. Modell der Dimensionsreduktion	51
4. Vorhersage in Linearen Modellen	55
4.1. Modelldefinition	55
4.2. Vorhersageintervalle	61
4.3. Weitere simultane Vorhersageintervalle unter Normalverteilungsannahme	64
5. Partielle kleinste Quadrate	69
5.1. Motivation	70
5.2. Modellspezifikation und Problemstellung	73
5.3. Die Vorhersage auf LP_1	75
5.3.1. Vorhersage durch Maximierung der Korrelation	75
5.3.2. Vorhersage durch Maximierung der Varianz (PCA-Methode)	75
5.3.3. Vorhersage durch Maximierung der Kovarianz: (PLS-Methode)	77
5.4. Modifizierte partielle kleinste Quadrate	78
5.5. Algorithmus der PLS-Methode: NIPALS-Algorithmus	80
5.6. Schlussbemerkung und Interpretation	82

6. Simulationsstudie	87
6.1. Modellspezifikation	87
6.2. Dimensionsreduktion	87
6.3. Simulationsdesign und Resultate	89
7. Ausblick	95
A. Lösungen einiger Maximierungsprobleme	97

Danksagung

Mein Dank gilt an dieser Stelle Herrn Prof. Schwabe, der diese Arbeit betreut und nach Fertigstellung begutachtet hat. Ebenso danke ich Herrn Prof. Schlittgen für die Erstellung des Gutachtens.

Den Professoren und Mitarbeitern des Instituts für Mathematische Stochastik danke ich für die Möglichkeit, im Rahmen des Oberseminars zu meinem Thema vortragen zu können. Verschiedene Hinweise und Anregungen, in diesem Rahmen haben ebenfalls zur Entstehung der Arbeit beigetragen.

Für eine schöne, abwechslungsreiche freundschaftliche Arbeitsatmosphäre, viele wertvolle Anregungen und stete Hilfsbereitschaft, die wesentlich zum Gelingen dieser Arbeit beigetragen haben, möchte ich mich den Kolleginnen und Kollegen bedanken.

Bedanken möchte ich mich auch bei Frau Altenkirch für alles. Der besten Freundin Rita danke ich aus ganzem Herzen für alles.

Weiterhin möchte ich den Menschen danken, die mir das alles letztendlichmöglich gemacht haben. Für die aufopferungsvolle Unterstützung auf den vielen Wegen, die ich bis zu diesem Punkt gehen musste, bedanke ich mich bei meinen Eltern, meinen Brüdern, meinen Schwestern und meinen Freunden.

Bei der Aleppo-Universität, bedanke ich mich für die finanzielle Unterstützung meines Studiums, ohne die diese Arbeit so nicht möglich gewesen wäre.

Zusammenfassung

Die Methode der partiellen kleinsten Quadrate (PLS = partial least squares) ist ein aktuelles Verfahren zur Dimensionsreduktion und Vorhersage in hochdimensionalen Datensätzen, wie sie z.B. in den Neurowissenschaften bei MRT-Daten zur Analyse von Hirnaktivitäten oder bei der Bildanalyse vorliegen. Es gibt eine Reihe sehr unterschiedlicher Algorithmen, die unter dem Begriff partial least squares zusammengefasst werden. Die statistischen Eigenschaften der Methode der partiellen kleinsten Quadrate sind jedoch bisher wenig erforscht. Dies liegt zu einem nicht unwesentlichen Teil an der Problematik, dass diese Methode der partiellen kleinsten Quadrate bisher nur unzureichend mit wahrscheinlichkeitstheoretischen Mitteln beschrieben werden kann.

In der vorliegenden Arbeit befassen wir uns mit einem Ansatz zur Entwicklung theoretischer Grundlagen der Methode der partiellen kleinsten Quadrate und wenden diese zur Vorhersage und zur Dimensionsreduktion an. Hierzu führen wir zuerst eine etwas allgemeinere Theorie der Vorhersage ein, wie sie in der Literatur bisher wenig behandelt wurde. Dazu geben wir zuerst einen Abriss der für die Arbeit benötigten Grundlagen der mathematischen Statistik. Dies beinhaltet insbesondere die Eigenschaften des besten erwartungstreuen Schätzers, d.h. desjenigen mit gleichmäßig kleinster Varianz. Danach stellen wir die in dieser Arbeit behandelten Assoziations- wesentlichen Eigenschaften der besten erwartungstreuen Vorhersage. Hierzu wird ein allgemeines Kriterium zur Ermittlung der besten Vorhersage abgeleitet. Im Anschluss werden einige aus der Literatur bekannte mathematische Charakterisierungen vorgestellt und durch neue ergänzt sowie ein neues analytisches Verfahren zur Verbesserung der Vorhersage entwickelt und auf seine Tauglichkeit überprüft. Anschließend wird mit Hilfe einer Simulationsstudie unter der Verteilungsannahme einer multivariaten Normalverteilung mit autoregressiver Kovarianzstruktur am Beispiel des Vorhersageintervalls die Qualität des neu entwickelten modifizierten Verfahrens mit dem Standardansatz der partiellen kleinsten Quadrate sowie mit Methoden der Hauptkomponentenanalyse und der kanonischen Korrelationsanalyse verglichen. Als Ergebnis demonstriert die Simulationsstudie die Vorteile des neu vorgeschlagenen modifizierten Ansatzes gegenüber den bekannten Methoden. Die abschließende Diskussion fasst die Ergebnisse zusammen und gibt einen Ausblick auf mögliche Erweiterungen des vorgestellten Konzepts.

Summary

Partial least squares is a topic method for dimension reduction and for prediction in high-dimensional data sets, for example in neuroscience with fMRI data for the analysis of brain activity. There are a number of different algorithms using the PLS approach. However, until now the statistical properties of the PLS method are not well investigated. This may be caused by the problem that the method of partial least squares is difficult to be described in terms of probability theory.

The aim of the present thesis is to develop a theoretical background of the Partial Least Squares methodology and to apply this approach to prediction and dimension reduction. In this context we also address a general theory of prediction, which has been neglected in the literature. We first give a summary of the background of mathematical statistics required for our purposes. In particular, properties of the uniformly minimum variance unbiased estimator are introduced. After that we present an introduction to the relevant association structures. Based on this we discuss the theory of prediction and, in particular, the properties of the best unbiased prediction. For this a general characterization for the best prediction is derived. Then we present and analyze some mathematical characterizations of the partial least squares approach given in the literature as well as some new ones and develop analytical solutions to improve the precision of the prediction. Subsequently we provide a simulation study, in which the performance of the newly proposed modified partial least squares method is compared with that of the standard least squares method as well as with the competing methods of principal component and canonical regression analysis under the model assumption of a multivariate normal distribution with an autoregressive covariance structure in the example of the construction of a prediction interval. As a result the simulation study demonstrates the advantages of the newly suggested modified approach compared to the methods existing so far in the literature. We conclude with a discussion of the results and an outlook to possible future work.

Einleitung

In vielen Anwendungsbereichen der Statistik wie der Medizin oder Ökonomie werden in zunehmendem Maße große und insbesondere hochdimensionale Datensätze erhoben, um damit interessierende Fragestellungen untersuchen zu können. Verfahren zur Dimensionsreduktion werden in der Statistik vielfältig verwendet:

1. Um die Daten auf geeignete Weise in einen Raum mit niedrigerer Dimension zu transformieren, in welchem die weitere Analyse und statistische Aufgaben (Punktschätzung, Intervallschätzung, Test, und Vorhersage) vorgenommen werden.
2. Um bessere Anpassung an die Daten und einfache Modelle zu erreichen.
3. Zum Einsparen von Speicherplatz
4. Zur Visualisierung und damit zum besseren Erkunden der Daten.

Die Regressionsanalyse stellt ein bekanntes Verfahren zur Dimensionsreduktion dar. Andere Methoden, die mit dem Begriff der Dimensionsreduktion verbunden werden, sind die Hauptkomponentenanalyse (engl. principal component analysis, PCA) und die kanonische Korrelation (engl. canonical correlation, CC). Ein aktuelleres Verfahren zur Dimensionsreduktion und Vorhersage in hochdimensionalen Datensätzen, wie sie z.B. in den Neurowissenschaften bei MRT-Daten zur Analyse von Hirnaktivitäten oder bei der Bildanalyse vorliegen, ist die Methode der partiellen kleinsten Quadrate (PLS = partial least squares). Diese Methode wurde ursprünglich von H. Wold (1966) eingeführt und wird in verschiedenen Anwendungsgebieten zur Regressions- oder Vorhersage für die Analyse hochdimensionaler Daten eingesetzt. Die statistischen Eigenschaften der Methode der partiellen kleinsten Quadrate sind jedoch bisher wenig erforscht. Dies liegt zu einem nicht unwesentlichen Teil an der Problematik, dass diese Methode der partiellen kleinsten Quadrate bisher nur unzureichend mit wahrscheinlichkeitstheoretischen Mitteln beschrieben werden kann. Ein mathematisches Modell des PLS-Ansatzes wurde von Helland (1988, 1990) eingeführt. Einige Erweiterungen des Resultates von Helland in mehrdimensionalen Modellen und Anwendungen für Klassifikation wurden von Liu und Rayens (2007) betrachtet. Die Shrinkage-Eigenschaft wurde von Krämer (2007) untersucht. Eine Erweiterung der PLS-Methode wurde von Li, Udén und v.Rosen (2013) modifiziert. In der vorliegenden Arbeit befassen wir uns mit einem Ansatz zur Entwicklung theoretischer Grundlagen der Methode der partiellen kleinsten Quadrate und wenden diese zur Vorhersage und zur Dimensionsreduktion an. Hierzu führen wir zuerst eine etwas allgemeinere Theorie der Vorhersage ein. Dazu geben wir zuerst einen Abriss der für die Arbeit benötigten Grundlagen der mathematischen Statistik. Dies

beinhaltet insbesondere die Eigenschaften des besten erwartungstreuen Schätzers, d.h. desjenigen mit gleichmäßig kleinster Varianz. Im zweiten Kapitel stellen wir die in dieser Arbeit behandelten Assoziationsstrukturen vor. Darauf aufbauend entwickeln wir im dritten Kapitel die Theorie der Vorhersage und leiten die wesentlichen Eigenschaften der besten erwartungstreuen Vorhersage her. Das vierte Kapitel behandelt die Vorhersage im linearen Modell. Im fünften Kapitel werden mathematische Charakterisierungen der PLS-Methode vorgestellt und durch neue ergänzt sowie ein neues analytisches Verfahren zur Verbesserung der Vorhersage entwickelt und auf seine Tauglichkeit überprüft. Anschließend wird mit Hilfe einer Simulationsstudie unter der Verteilungsannahme einer multivariaten Normalverteilung mit autoregressiver Kovarianzstruktur am Beispiel des Vorhersageintervalls die Qualität des neu entwickelten modifizierten Verfahrens mit dem Standardansatz der partiellen kleinsten Quadrate sowie mit Methoden der Hauptkomponentenanalyse und der kanonischen Korrelationsanalyse verglichen. Als Ergebnis demonstriert die Simulationsstudie die Vorteile des neu vorgeschlagenen modifizierten Ansatzes gegenüber den bekannten Methoden. Zusammenfassend kann gesagt werden, dass diese Arbeit das Ziel verfolgt, und diese Interpretation bietet mögliche Erweiterungen der vorgestellten Konzepte für verschiedene Modelle zu entwickeln.

1. Allgemeine Grundlagen der mathematischen Statistik

In diesem Kapitel stellen wir die wichtigsten mathematischen Begriffe vor, auf denen die Entwicklungen in den folgenden Kapiteln beruhen. Dazu befassen wir uns mit Theorie und Anwendung der statistischen Modelle. Schließlich geben wir einen kurzen Überblick über die wichtigsten Resultate zur Parameterschätzung im linearen Normalverteilungsmodell und im Modell mit multivariaten, u.i.v. normalverteilten ZV'en und verweisen für Beweise und weitere Ergebnisse auf die Literatur, z.B. Shao (2003), Alvin (2008), Christensen (2011), Christensen und Lin (2013) und Hocking (2005). Zuletzt beleuchten wir kurz das AR(1)-Modell.

1.1. Statistisches Experiment

Grundlage der mathematischen Behandlung einer statistischen Problemstellung bildet das statistische Experiment welches das den Beobachtungen zugrundegelegte Modelle beschreibt.

Definition 1.1 (*statistisches Experiment*)

Ein statistisches Experiment ist ein Tripel $E = (M, \mathcal{B}, (P_\theta)_{\theta \in \Theta})$ bestehend aus

- einer nichtleeren Menge M
- einer σ -Algebra \mathcal{B} über M , die die beobachtbaren Ereignisse enthält .
- einer Familie $(P_\theta)_{\theta \in \Theta}$ von Wahrscheinlichkeitsmaßen auf (M, \mathcal{B}) , die mit Element des Parameterraums Θ parametrisiert ist. Mit jedem statistischen Experiment kann die Beobachtung einer Zufallsvariablen

$$X : (\Omega, \mathcal{C}) \rightarrow (M, \mathcal{B})$$

und eine Familie $(Q_\theta)_{\theta \in \Theta}$ von W -Maßen auf (Ω, \mathcal{C}) verbunden werden, so dass

$$Q_\theta^X = P_\theta \quad \text{für alle } \theta \in \Theta.$$

Die Formulierung von statistischen Experimenten wird anhand der folgenden Beispiele illustriert.

Beispiele für statistische Experimente

Das statistische Experiment beschreibt stets das Ergebnis eines Zufallsexperiments, etwa die Werte einer erhaltenen Stichprobe oder gesammelte Messergebnisse eines Experiments. Somit ist die Verteilung der Zufallsvariable das Schlüsselement. Das statistische Experiment ist dann eine geeignete Familie von solchen Verteilungen. Anhand von drei Beispielen wird im Folgenden die Formulierung von statistischen Experimenten illustriert.

(1) Vereinfachtes Normalverteilungsmodell

Wir betrachten das Modell mit $n \geq 2$ u.i.v. $N(\beta, \sigma^2)$ -verteilten Zufallsvariablen X_1, X_2, \dots, X_n , wobei $\beta \in \mathbb{R}$ und $\sigma^2 \in (0, \infty)$ die Parameter sind. Mit \mathbf{X} bezeichnen wir die \mathbb{R}^n -wertige Zufallsvariable $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$. Damit hat man das statistische Experiment

$$E = \left(\mathbb{R}^n, \mathcal{B}^n, \left(\otimes_{i=1}^n N(\beta, \sigma^2) \right)_{(\beta, \sigma^2) \in \mathbb{R} \times (0, \infty)} \right) \quad (1.1)$$

(2) Lineares Normalverteilungsmodell

Im Folgenden sei ein lineares Normalverteilungsmodell zu Grund gelegt:

$$\mathbf{X} \sim (\mathbf{B}\boldsymbol{\beta}, \sigma^2 \mathbf{V}), \quad \mathbf{V} \text{ bekannt}$$

Die Zufallsvariable $\mathbf{X} = (X_1, X_2, \dots, X_n)^T$ des Modells ist \mathbb{R}^n -wertig, wobei \mathbf{B} eine gegebene reelle $n \times p$ Matrix vom Rang p , und $n > p$ vorausgesetzt ist, mit $\mathbf{V} \in PD(n)$ (Die Menge aller positiv definiten reellen $(n \times n)$ -Matrizen wird mit $PD(n)$ abgekürzt). Der Parameter ist $\boldsymbol{\theta} = (\boldsymbol{\beta}^T, \sigma^2)^T \in \mathbb{R}^p \times (0, \infty)$

$$E = \left(\mathbb{R}^n, \mathcal{B}^n, N_n(\mathbf{B}\boldsymbol{\beta}, \sigma^2)_{(\boldsymbol{\beta}, \sigma^2) \in \mathbb{R}^p \times (0, \infty)} \right) \quad (1.2)$$

(3) Modell mit multivariaten, u.i.v. normalverteilten ZV'en

Seien X_1, X_2, \dots, X_n stochastisch unabhängig und identisch $N_p(\boldsymbol{\beta}, \mathbf{V})$ -verteilte \mathbb{R}^p -wertige Zufallsvariablen, wobei $(\boldsymbol{\beta}, \mathbf{V}) \in \mathbb{R}^p \times PD(p)$ die Parameter sind. Kurz:

$$E = \left((\mathbb{R}^p)^n, (\mathcal{B}^p)^n, \left(\otimes_{i=1}^n N(\boldsymbol{\beta}, \mathbf{V}) \right)_{(\boldsymbol{\beta}, \mathbf{V}) \in \mathbb{R}^p \times PD(p)} \right) \quad (1.3)$$

Um die Aufgabenstellungen weiter zu erläutern und nachfolgende Begriffsbildungen vorzunehmen, werden wir im weiteren Verlauf unserer Betrachtungen neben dem statistischen Experiment einen weiteren Begriff benötigen, und zwar den der statistischen Entscheidungsprobleme

Definition 1.2 (*statistisches Entscheidungsproblem*)

Ein statistisches Entscheidungsproblem ist ein Tripel (E, A, L) bestehend aus

- einen statistischen Experiment $E = (M, \mathcal{B}, (P_\theta)_{\theta \in \Theta})$

- einem Antwortraum A
- einer Verlustfunktion $L : A \times \Theta \rightarrow [0, \infty)$

Falls die Parametermenge eine Teilmenge des \mathbb{R}^d ist, spricht man von einem parametrischen, ansonsten von einem nichtparametrischen Modell. Die Zustandsräume

$$\mathbb{R}^p \times (0, \infty), \quad \mathbb{R}^p \times PD(p)$$

implizieren zum Beispiel parametrische Modelle. In dieser Arbeit beschränken wir uns im Wesentlichen auf parametrische Modelle.

1.2. Schätzproblem

Wir gehen wieder von einem allgemeinen statistischen Experiment aus,

$$E = \left(M, \mathcal{B}, (P_\theta^X)_{\theta \in \Theta} \right)$$

desweiteren sei ein Aspekt des Parameters gegeben, der einen interessierenden Teilparameter $g(\theta)$ beschreibt, $g : \Theta \rightarrow A$, wir nennen den Aspekt $g(\theta), \theta \in \Theta$ ebenfalls einen Parameter. Unter einem Schätzer für $g(\theta)$ verstehen wir eine Funktion

$$T : (M, \mathcal{B}) \rightarrow (A, \mathcal{A}),$$

die also für jeden möglichen Wert $x \in M$ der Zufallsvariablen X die Schätzung $T(x)$ für $g(\theta)$ liefert, und T messbar ist. Sei

$$D = \{T | T : M \rightarrow A \text{ messbar}, E(T^2(X)) < \infty\} \quad (1.4)$$

die Menge aller Schätzer. Um den stochastischen Kontext zu betonen, unterscheiden wir:

- $T(X)$ Schätzer für $g(\theta)$ Zufallsvariable auf (M, \mathcal{A}) (und die im Wesentlichen dasselbe wie T ist).
- $T(x)$ Schätzung für $g(\theta)$ Realisation.

Mit der Bezeichnung T statt \hat{g} eines Schätzers für $g(\theta)$ haben wir einen kleinen Notationswechsel vorgenommen, der natürlich inhaltlich Bedeutung hat. In dieser Arbeit betrachten wir (und vergleichen wir) mehrere/viele Schätzer für denselben Parameter $g(\theta)$, und Bezeichnungen wie T, T_1, T^* etc, für verschiedene Schätzer für $g(\theta)$. Das Ziel ist es, die Qualität eines Schätzers $T(X)$ für den Parameter $g(\theta)$ zu messen. Als Messung des Fehlers stellen wir zunächst die Verlustfunktion im ein- mehrdimensionale Fall vor und behandeln danach Risikofunktion. Abschließend wird sein, die Qualität des Schätzers anhand seiner Streuung zu messen.

Einige Verlustfunktionen:

(i) Der eindimensionale Fall ($A = \mathbb{R}$):

$$L(z, \theta) = |z - g(\theta)| \quad (\text{absolute Abweichung}),$$

$$L(z, \theta) = (z - g(\theta))^2 \quad (\text{quadratische Abweichung}),$$

$$L(z, \theta) = \frac{(z - g(\theta))^2}{g(\theta)^2} \quad (\text{relative quadratische Abweichung}),$$

$$L(z, \theta) = (w(\theta))^2 (z - g(\theta))^2 \quad \text{gewichtete quadratische Abweichung, mit einer gegebenen Gewichtsfunktion } w : \Theta \rightarrow [0, \infty).$$

(ii) Der mehrdimensionale Fall ($A = \mathbb{R}^p$):

$$L(\mathbf{z}, \boldsymbol{\theta}) = (\mathbf{z} - \mathbf{g}(\boldsymbol{\theta}))^T (\mathbf{z} - \mathbf{g}(\boldsymbol{\theta})) \quad (\text{quadratische Abweichung}),$$

$$L(\mathbf{z}, \boldsymbol{\theta}) = \frac{(\mathbf{z} - \mathbf{g}(\boldsymbol{\theta}))^T (\mathbf{z} - \mathbf{g}(\boldsymbol{\theta}))}{\mathbf{g}(\boldsymbol{\theta})^T \mathbf{g}(\boldsymbol{\theta})} \quad (\text{relative quadratische Abweichung}),$$

$$L(\mathbf{z}, \boldsymbol{\theta}) = (\mathbf{z} - \mathbf{g}(\boldsymbol{\theta}))^T \mathbf{W}(\boldsymbol{\theta}) (\mathbf{z} - \mathbf{g}(\boldsymbol{\theta})) \quad \text{gewichtete quadratische Abweichung, mit einer gegebenen Gewichtstransformation } \mathbf{W} : \Theta \rightarrow PD(p).$$

1.2.1. Optimalitätskriterien

Zum Vergleich von Schätzer interessieren wir uns oft für die Frage: wie soll man das Verhalten eines Schätzers $T(X)$ für $g(\theta)$ messen? Daher gehen wir von einer gegebenen Verlustfunktion L aus

$$L : A \times \Theta \rightarrow [0, \infty)$$

Der Vergleich verschiedener Schätzer für $g(\theta)$ orientiert sich an deren Risikofunktion.

Definition 1.3 (Risikofunktion)

Sei (E, A, L) ein statistisches Entscheidungsproblem. Dann heißt

$$R(T(X), g(\theta)) = E_{\theta}(L(T(X), g(\theta))) \tag{1.5}$$

die Erwartung von einer Verlustfunktion des Schätzers $T(X)$ die Risikofunktion des Schätzers $T(X)$.

Das Ziel wird sein, ein Kriterium zu finden, welches bereits vor der Datenerhebung zur Beurteilung eines Schätzers genutzt werden kann. Dafür kommt unter anderem das im Folgenden vorgestellte Maß des mittleren quadratischen Fehlers in Frage.

Der mittlere quadratische Fehler, kurz MSE (engl. mean squared error), ist ein Gütemaß für den Punktschätzer. Er setzt sich zusammen aus dem Bias und der Varianz des Punktschätzers. Hauptsächlich betrachten wir in diesem Abschnitt den quadratischen Verlust.

Definition 1.4 (mittlerer quadratischer Fehler)

Der mittlere quadratische Fehler eines Punktschätzers $T \in D$ für einen Parameter $g(\theta)$, ist definiert als

$$MSE_{\theta}(T(X), g(\theta)) = \mathbb{R}(T(X), g(\theta)) = E_{\theta}[(T(X) - g(\theta))^2] \tag{1.6}$$

Weiterhin heißt

$$b_{\theta}(T(X), g(\theta)) = E_{\theta}(T(X)) - g(\theta)$$

die Bias von T bzgl g . Gilt $b_{\theta}(T(X), g(\theta)) = 0, \forall \theta \in \Theta$, wenn T unverzerrt ist.

Daraus als Korollar zur Definition 1.4 ergibt sich die folgende wichtige Zerlegung des mittleren quadratischen Fehlers in Varianz des Schätzers und Quadrat des Biases:

$$\begin{aligned} MSE_{\theta}(T(X), g(\theta)) &= Var_{\theta}(T(\theta)) + [E_{\theta}(T(X)) - g(\theta)]^2 \\ &= Var_{\theta}(T(X)) + b_{\theta}^2(T(X), g(\theta)). \end{aligned}$$

Interessant in der statistischen Theorie ist der mehrdimensionale Fall, daher betrachten wir die Verallgemeinerungen der mittleren quadratischen Abweichung. Die Zweckmäßigkeit von verzerrten Schätzer im mehrparametrischen Fall lässt sich aus dem MSE-Kriterium herleiten. Es ist definiert durch:

Definition 1.5 (*R^p -wertiger Parameter und MSE*) Schlittgen (2009, S.98,99)

Im Fall $\mathbf{g}(\boldsymbol{\theta}) \in \mathbb{R}^p$ nennt man

(a) den Ausdruck

$$MSEM_{\boldsymbol{\theta}} = E_{\boldsymbol{\theta}}[(\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta}))(\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta}))^T]$$

die Matrix des mittleren quadratischen Fehlers.

(b)

$$MSE_{\boldsymbol{\theta}} = E_{\boldsymbol{\theta}}(\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta}))^T (\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta})) \quad (1.7)$$

den mittleren quadratischen Fehler

Beachte:

- Die Matrix des mittleren quadratischen Fehlers hat die einzelnen MSE's auf der Hauptdiagonalen.
- Er lässt sich durch Ergänzen und Ausmultiplizieren einfach umformen zu

$$\begin{aligned} MSE_{\boldsymbol{\theta}} &= E_{\boldsymbol{\theta}}[(\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta}))^T ((\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta})))] \\ &= Spur(E_{\boldsymbol{\theta}}[(\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta}))((\mathbf{T}(\mathbf{X}) - \mathbf{g}(\boldsymbol{\theta}))^T)]) \\ &= \sum_{i=1}^p MSE_{\boldsymbol{\theta}}(T_i(\mathbf{X}), g_i(\boldsymbol{\theta})) \quad (T_i, g_i \text{ der } i\text{-ten Komponente}) \\ &= \sum_{i=1}^p Var(T_i(\mathbf{X})) + (E[(\mathbf{T}(\mathbf{X}))] - \mathbf{g}(\boldsymbol{\theta}))^T (E[(\mathbf{T}(\mathbf{X}))] - \mathbf{g}(\boldsymbol{\theta})) \end{aligned}$$

1.2.2. Eigenschaften des UMVUE-Schätzers

UMVUE steht für engl. Uniformly Minimum Variance Unbiased Estimator. Für eine vorgegebene Funktion g des Parameters θ betrachten wir die Menge aller erwartungstreuen Schätzer für $g(\theta)$ mit endlicher Varianz

$$U_g = \{T|T : M \rightarrow A \text{ messbar, } E_{\theta}(T(X)) = g(\theta), E_{\theta}(T^2(X)) < \infty\}$$

1. Allgemeine Grundlagen der mathematischen Statistik

und die Menge aller Schätzer für $g(\theta)$ mit endlicher Varianz und Erwartungswert gleich null

$$U_0 = \{T|T : M \rightarrow A \text{ messbar, } E_\theta(T(X)) = 0, E_\theta(T^2(X)) < \infty\}.$$

Definition 1.6 (UMVUE-Schätzer)

Eine Statistik $T^* \in U_g$ heißt UMVUE, wenn gilt

$$\text{Var}(T^*) = \inf\{\text{Var}(T) : T \in U_g\}.$$

- (a) Der Satz von Rao-Blackwell sagt aus, dass durch Bildung der bedingten Erwartung bezüglich einer suffizienten Statistik eines Schätzers mit gleichmäßig nicht größerem Risiko entsteht.

Theorem 1.1 (Rao-Blackwell-Verbesserung) Witting (1985, S. 349)

Seien $T(X)$ eine suffiziente Statistik und $S(X)$ ein Schätzer mit $S \in U_g$ (S im Wesentlichen ist dasselbe wie $S(X)$). Definiere die Rao-Blackwell-Verbesserung $S^*(X) = E[S(X)|T(X)]$. Dann gilt:

1. $S^* \in U_g$
2. $\text{Var}_\theta(S^*(X)) \leq \text{Var}_\theta(S(X)), \forall \theta \in \Theta$, d.h. S^* ist besser als S .

- (b) Wenn die Statistik $T(X)$ suffizient sogar vollständig ist, so für $S \in U_g$ liefert die Rao-Blackwell-Verbesserung $S^*(X) = E[S(X)|T(X)]$ bereits ein UMVUE-Schätzer, wie der nächste Satz von Lehmann-Scheffé aussagt:

Theorem 1.2 (Lehmann-Scheffé-Verbesserung) Witting (1985, S. 354)

Sei $T(X)$ vollständige suffiziente Statistik und $S \in U_g$. Dann ist der Schätzer

$$S^*(X) = E[S(X)|T(X)]$$

gleichmäßig optimal in der Menge aller erwartungstreuen Schätzer für $g(\theta)$, d.h. $S^*(X)$ ist ein erwartungstreuer Schätzer für $g(\theta)$ mit gleichmäßig minimaler Varianz (UMVUE).

Der Satz von Lehmann-Scheffé liefern ein Verfahren, einen gleichmäßig besten erwartungstreuen Schätzer zu finden. Mittels der Fisher-Information $I(\theta)$ ergibt sich nun eine untere Schranke für die Varianz eines Schätzers $T^* \in U_g$.

- (c) **Theorem 1.3 (Fisher-Information und UMVUE) Shao (2003 S. 171,172)**

Falls $\mathcal{P} = \{f_\theta, \theta \in \Theta\}$ eine exponentielle Familie mit Darstellung

$$f_\theta(x) = \exp(c(\theta)T(x) + d(\theta) + l(x))$$

ist, und $T^*(X)$ ein erwartungstreuer Schätzer für $g(\theta)$, so dass

$$\text{Var}_\theta(T^*(X)) = \frac{(\frac{\partial}{\partial \theta} g(\theta))^2}{I(\theta)}, \forall \theta \in \Theta. \quad (1.8)$$

Dann

1. $T^*(x)$ UMVUE-Schätzung für $g(\theta)$.
2. $T^*(x)$ lineare Funktion in $T(x)$

Die Optimalität der erwartungstreuen Schätzer lässt sich mit der Kovarianzmethode charakterisieren

(d) Theorem 1.4 (Kovarianzmethode) Shao (2003, Kap. 3.1.2)

Sei $T^* \in U_g$ erwartungstreuer Schätzer, dann gilt:

$$\text{Cov}_\theta(T^*(X), T_0(X)) = 0, \forall \theta \in \Theta \text{ und } \forall T_0 \in U_0 \Leftrightarrow T^*(X) \text{ UMVUE für } g(\theta) \quad (1.9)$$

Aus Theorem 1.4 ergibt sich insbesondere die folgende Folgerung.

Folgerung 1.1 Sei $T^* \in U_g$ ein erwartungstreuer Schätzer, dann gilt

$$E_\theta(T^*(X)T_0(X)) = 0, \forall T_0 \in U_0 \text{ und } \forall \theta \in \Theta \Leftrightarrow T^*(X) \text{ UMVUE.} \quad (1.10)$$

Beachte:

- Wenn T^* UMVUE-Schätzer ist, gilt

$$\text{Var}_\theta(T^*) = \text{Cov}_\theta(T^*, T), \forall T \in U_g.$$

und die Varianz von $T - T^*$ ist gleich die Differenz der Varianzen

$$\begin{aligned} \text{Var}_\theta(T - T^*) &= \text{Var}_\theta(T^*) + \text{Var}_\theta(T) - 2 \text{Cov}_\theta(T, T^*) \\ &= \text{Var}_\theta(T) - \text{Cov}_\theta(T^*, T). \\ &= \text{Var}_\theta(T) - \text{Var}_\theta(T^*). \end{aligned}$$

- Hieraus ergibt sich insbesondere

$$\frac{\text{Var}_\theta(T - T^*)}{\text{Var}_\theta(T)} = 1 - \text{Corr}_\theta^2(T, T^*).$$

Mit der Schreibweise

$$\text{Corr}_\theta^2(T, T^*) = 1 - \frac{\text{Var}_\theta(T - T^*)}{\text{Var}_\theta(T)}, \text{ in dem } 0 \leq \text{Corr}_\theta^2(T, T^*) \leq 1.$$

ergibt sich:

$$\text{Corr}_\theta^2(T, T^*) = 1 \iff \text{Var}_\theta(T - T^*) = 0 \text{ d.h. } T = T^* \text{ } P_\theta \text{ fast sicher}$$

- Wenn T_1, T_2 zwei UMVUE-Schätzer für den Parameter $g(\theta)$ sind, dann ist $T_1 = T_2$ P_θ fast sicher, da

$$\begin{aligned} E_\theta[(T_1 - T_2)^2] &= E_\theta[T_1(T_1 - T_2)] + E_\theta[T_2(T_2 - T_1)] \\ &= 0 \quad (\text{Folgerung(1.1)}) \end{aligned}$$

Beispiel 1.1 Vereinfachtes Normalverteilungsmodell: Betrachte das Modell

$$\begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim N_2 \left(\begin{pmatrix} 1 \\ 1 \end{pmatrix} \beta, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right), \text{ mit } \mathbf{V} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \text{ bekannt und } |\rho| < 1.$$

Parameter des Modells ist: $\beta \in \mathbb{R}$.

- Der gewichtete kleinste Quadrate Schätzer (WLS), bzw. ML-Schätzer ist:

$$\begin{aligned} \hat{\beta}_{WLS} &= \arg \min_{\beta} (\mathbf{X} - \mathbf{1}_2 \beta)^T \mathbf{V}^{-1} (\mathbf{X} - \mathbf{1}_2 \beta) \\ &= (\mathbf{1}_2^T \mathbf{V}^{-1} \mathbf{1}_2)^{-1} \mathbf{1}_2^T \mathbf{V}^{-1} \mathbf{X} \\ &= \left(\mathbf{1}_2^T \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix} \mathbf{1}_2 \right)^{-1} \mathbf{1}_2^T \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix} \mathbf{X} \\ &= \frac{1}{2(1-\rho)} (1-\rho)(X_1 + X_2) \\ &= \bar{X} \end{aligned}$$

- Der gewöhnliche kleinste Quadrate Schätzer (OLS) ist:

$$\begin{aligned} \hat{\beta}_{OLS} &= \arg \min_{\beta} (\mathbf{X} - \mathbf{1}_2 \beta)^T (\mathbf{X} - \mathbf{1}_2 \beta) \\ &= (\mathbf{1}_2^T \mathbf{1}_2)^{-1} \mathbf{1}_2^T \mathbf{X} \\ &= \frac{(X_1 + X_2)}{2} \\ &= \bar{X} \end{aligned}$$

Interpretation im Fall $A = \mathbb{R}^p$

Im Fall $A = \mathbb{R}^p$ steht MCUE für Minimum Covariance unbiased estimator. Betrachtet man nur unverzerrte Schätzer, so kann man die Kovarianz des Schätzers als Maß für die Qualität des Schätzers heranziehen, da unter Unverzerrtheit die Kovarianz des Schätzers gleich der mittleren quadratischen Fehler ist.

Definition 1.7 (Kovarianz-Matrix des Schätzers) Rao, Shalabh, Toutenburg und Heumann (2008, Kap.3.4)

Sei $\mathbf{T} = (T_1, \dots, T_p)^T$ ein unverzerrter Schätzer für den \mathbb{R}^p -Parameter

$$\mathbf{g}(\boldsymbol{\theta}) = (g_1(\boldsymbol{\theta}), \dots, g_p(\boldsymbol{\theta}))^T.$$

Dann heißt die $p \times p$ -Matrix $\text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}) = E_{\boldsymbol{\theta}}[(\mathbf{T} - \mathbf{g}(\boldsymbol{\theta}))(\mathbf{T} - \mathbf{g}(\boldsymbol{\theta}))^T]$ Kovarianz-Matrix von \mathbf{T} .

Ein Schätzer ist in diesem Sinn besser als alle anderen unverzerrten Schätzer, falls seine Kovarianz minimal ist, was zu folgender Optimalitätseigenschaft führt.

Definition 1.8 (MCUE-Schätzer) Rao, Shalabh, Toutenburg und Heumann (2008, Kap.3.4)

Ein Schätzer $\mathbf{T} \in U_g$ von \mathbb{R}^p -Parameter $\mathbf{g}(\boldsymbol{\theta})$ heißt MCUE, wenn

$$\text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}) \leq_{\text{psd}} \text{Cov}_{\boldsymbol{\theta}}(\mathbf{S}), \quad \forall \boldsymbol{\theta} \in \Theta \quad \text{und} \quad \forall \mathbf{S} \in U_g. \quad (1.11)$$

gilt, wobei \leq_{psd} die Anordnung im Sinne der positiven Definitheit ist.

Lemma 1.1 Rao, Shalabh, Toutenburg und Heumann (2008, Kap.3.4)

Die folgenden Bedingungen sind äquivalent:

1. $\mathbf{T}^* = (T_1^*, \dots, T_p^*)^T$ ein MCUE-Schätzer für den \mathbb{R}^p -Parameter

$$\mathbf{g}(\boldsymbol{\theta}) = (g_1(\boldsymbol{\theta}), \dots, g_p(\boldsymbol{\theta}))^T$$

2. $\mathbf{c}^T \mathbf{T}^*$ ist UMVUE-Schätzer für den Parameter $\mathbf{c}^T \mathbf{g}(\boldsymbol{\theta}), \forall \mathbf{c} \in \mathbb{R}^p$.

3. $\text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}_0, \mathbf{T}^*) = \mathbf{0}_{q \times q}, \forall \mathbf{T}_0 \in U_0$

Bemerkung 1.1 Wenn \mathbf{T}^* MCUE-Schätzer ist, dann gilt:

$$\text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}^*) = \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}^*, \mathbf{T}), \quad \forall \mathbf{T} \in U_g$$

$$\begin{aligned} \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T} - \mathbf{T}^*) &= \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}^*) + \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}) - 2 \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}, \mathbf{T}^*) \\ &= \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}) - \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}^*, \mathbf{T}). \\ &= \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}) - \text{Cov}_{\boldsymbol{\theta}}(\mathbf{T}^*) \geq 0. \end{aligned}$$

Bemerkung 1.2 (BLUE-Schätzer)

Ein linearer UMVUE-Schätzer (MCUE-Schätzer) heißt BLUE-Schätzer (best linear unbiased estimator)

1.3. Verteilungsmodelle

In diesem Abschnitt stellen wir \mathbb{R}^{p+q} -wertige verteilte Zufallsvariablen (ZV) vor.

(a) Seien $\mathbf{X} : (M, \mathcal{A}) \longrightarrow (\mathbb{R}^p, \mathcal{B}^p), \mathbf{Y} : (M, \mathcal{A}) \longrightarrow (\mathbb{R}^q, \mathcal{B}^q),$

$$\begin{aligned} &\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \text{ eine } \mathbb{R}^{p+q} \text{ - wertige ZV, mit} \\ \boldsymbol{\mu} &= E \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix} \in \mathbb{R}^{p+q} \text{ bekannt} \\ \mathbf{V} &= \text{Cov} \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \in PD(p+q) \text{ bekannt.} \end{aligned}$$

(b) \mathbf{X} und \mathbf{Y} gemeinsam normalverteilt:

$$\begin{aligned} \mathbf{X} &\sim N_p(\boldsymbol{\mu}_x, \mathbf{V}_{xx}) \text{ und } \mathbf{Y} \sim N_q(\boldsymbol{\mu}_y, \mathbf{V}_{yy}) \\ \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} &\sim N_{p+q}(\boldsymbol{\mu}, \mathbf{V}), \text{ mit } \boldsymbol{\mu} \in \mathbb{R}^{p+q} \text{ und } \mathbf{V} \in PD(p+q) \end{aligned}$$

Wir partitionieren $\boldsymbol{\mu}$ und \mathbf{V} entsprechend der Teil-Dimensionen p und q , wie im Teil (a).

Der bedingte Erwartungswert von \mathbf{Y} unter \mathbf{X} ist gegeben durch:

$$E(\mathbf{Y}|\mathbf{X}) = \boldsymbol{\mu}_y + \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}(\mathbf{X} - \boldsymbol{\mu}_x). \quad (1.12)$$

Die bedingte Verteilung von \mathbf{Y} unter \mathbf{X} ist gegeben durch:

$$\mathbb{P}^{\mathbf{Y}|\mathbf{X}=\mathbf{x}} = N(\boldsymbol{\mu}_y + \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}(\mathbf{x} - \boldsymbol{\mu}_x), \mathbf{V}_{yy} - \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}\mathbf{V}_{yx}^T), \quad \mathbf{x} \in \mathbb{R}^p \quad (1.13)$$

1.4. Lineares Normalverteilungsmodell

Sei ein lineares Normalverteilungsmodell zu Grunde gelegt

$$\begin{aligned} \mathbf{X} &\sim N_n(\mathbf{B}\boldsymbol{\beta}, \sigma^2 \mathbf{V}), \quad (\boldsymbol{\beta}, \sigma^2) \in \mathbb{R}^p \times (0, \infty) \\ \mathbf{B} &\text{ eine feste } n \times p \text{ - Matrix, } \text{Rang}(\mathbf{B}) = p \\ \mathbf{V} &\in PD(n) \text{ bekannt, mit } p < n. \end{aligned}$$

Theorem 1.5 Als ML-Schätzung $(\widehat{\boldsymbol{\beta}}_{ML}(\mathbf{x}), \widehat{\sigma}_{ML}^2(\mathbf{x}))$ für $(\boldsymbol{\beta}, \sigma^2)$ ergibt sich (wobei wir $n > p$ voraussetzen):

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_{ML}(\mathbf{x}) &= (\mathbf{B}^T\mathbf{V}^{-1}\mathbf{B})^{-1}\mathbf{B}^T\mathbf{V}^{-1}\mathbf{x}, \\ \widehat{\sigma}_{ML}^2(\mathbf{x}) &= \frac{1}{n}(\mathbf{x} - \mathbf{B}(\mathbf{B}^T\mathbf{V}^{-1}\mathbf{B})^{-1}\mathbf{B}^T\mathbf{V}^{-1}\mathbf{x})^T(\mathbf{x} - \mathbf{B}(\mathbf{B}^T\mathbf{V}^{-1}\mathbf{B})^{-1}\mathbf{B}^T\mathbf{V}^{-1}\mathbf{x}). \end{aligned}$$

- Als Gauß-Markov-Schätzung $\widehat{\boldsymbol{\beta}}_{GM}(\mathbf{x})$ für $\boldsymbol{\beta}$ ergibt sich (wobei \mathbf{V} bekannt ist):

$$\widehat{\boldsymbol{\beta}}_{GM}(\mathbf{x}) = (\mathbf{B}^T\mathbf{V}^{-1}\mathbf{B})^{-1}\mathbf{B}^T\mathbf{V}^{-1}\mathbf{x} \quad (1.14)$$

$$\text{BLUE-Schätzung für } \boldsymbol{\beta}. \quad (1.15)$$

- Als gewöhnliche kleinste-Quadrate-Schätzung $\widehat{\boldsymbol{\beta}}_{KQ}(\mathbf{x})$ für $\boldsymbol{\beta}$ ergibt sich:

$$\widehat{\boldsymbol{\beta}}_{KQ}(\mathbf{x}) = (\mathbf{B}^T\mathbf{B})^{-1}\mathbf{B}^T\mathbf{x}. \quad (1.16)$$

- Es gilt auch :

$$\begin{pmatrix} \widehat{\beta}_{GM} \\ \widehat{\beta}_{KQ} \end{pmatrix} \sim N_{2p} \left(\begin{pmatrix} \beta \\ \beta \end{pmatrix}, \sigma^2 \begin{pmatrix} (\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} & (\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} \\ (\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} & (\mathbf{B}^T \mathbf{B})^{-1} (\mathbf{B}^T \mathbf{V} \mathbf{B}) (\mathbf{B}^T \mathbf{B})^{-1} \end{pmatrix} \right)$$

Hieraus ergibt sich auch:

$$\begin{aligned} \text{Cov}(\widehat{\beta}_{GM}, \widehat{\beta}_{GM} - \widehat{\beta}_{KQ}) &= \text{Cov}(\widehat{\beta}_{GM}) - \text{Cov}(\widehat{\beta}_{GM}, \widehat{\beta}_{KQ}) \\ &= (\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} - (\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} \\ &= \mathbf{0} \end{aligned}$$

Damit ist die Behauptung von Theorem(1.3) erfüllt.

Theorem 1.6 1) Als BLUE-Schätzung $\widehat{\beta}(\mathbf{x})$ für β ergibt sich:

$$\begin{aligned} \widehat{\beta}(\mathbf{x}) &= (\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{V}^{-1} \mathbf{x} \\ &= \text{MCUE-Schätzung für } \beta. \end{aligned}$$

2) Als UMVUE-Schätzung $\widehat{\sigma}^2(\mathbf{x})$ für σ^2 ergibt sich:

$$\widehat{\sigma}^2(\mathbf{x}) = \frac{1}{n-p} (\mathbf{x} - \mathbf{B}(\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{V}^{-1} \mathbf{x})^T (\mathbf{x} - \mathbf{B}(\mathbf{B}^T \mathbf{V}^{-1} \mathbf{B})^{-1} \mathbf{B}^T \mathbf{V}^{-1} \mathbf{x})$$

1.5. Modell mit multivariaten, u.i.v. normalverteilten ZV'en

Seien $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ stochastisch unabhängig und identisch $N_p(\beta, \mathbf{V})$ -verteilte \mathbb{R}^p -wertige Zufallsvariablen, wobei $(\beta, \mathbf{V}) \in \mathbb{R}^p \times PD(p)$ die Parameter sind. Wir fassen die Zufallsvariablen zu einer $\mathbb{R}^{n \times p}$ -wertigen Zufallsvariablen zusammen,

$$\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^T.$$

Die Lebesgue-Dichte der Verteilung $P_{\beta, \mathbf{V}}^{\mathbf{X}}$ ist gegeben durch:

$$\begin{aligned} f_{\beta, \mathbf{V}}(\mathbf{x}) &= \prod_{i=1}^n \left[(2\pi)^{-p/2} \det(\mathbf{V})^{-1/2} \exp \left(-\frac{1}{2} (\mathbf{x}_i - \beta)^T \mathbf{V}^{-1} (\mathbf{x}_i - \beta) \right) \right] \\ &= (2\pi)^{-np/2} \det(\mathbf{V})^{-n/2} \exp \left(-\frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \beta)^T \mathbf{V}^{-1} (\mathbf{x}_i - \beta) \right) \end{aligned}$$

für alle $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}^{n \times p}$, wobei wir $n > p$ voraussetzen.

Die Log-Likelihood Funktion ist

$$l_{\mathbf{x}}(\beta, \mathbf{V}) = -\frac{np}{2} \ln(2\pi) - \frac{n}{2} \ln(\det(\mathbf{V})) - \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i - \beta)^T \mathbf{V}^{-1} (\mathbf{x}_i - \beta).$$

1. Allgemeine Grundlagen der mathematischen Statistik

Wir betrachten die Statistik

$$\begin{aligned}\bar{\mathbf{x}} &= \frac{1}{n} \mathbf{x}^T \mathbf{1}_n \\ &= \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \\ \mathbf{S}(\mathbf{x}) &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \\ &= \mathbf{x}^T \mathbf{x} - \bar{\mathbf{x}} \bar{\mathbf{x}}^T.\end{aligned}$$

Elementare Eigenschaften von \mathbf{S} und $\bar{\mathbf{x}}$ Schlittgen (2009, S. 101)

1. Es lässt sich zeigen, dass $\bar{\mathbf{X}} \sim N_p(\boldsymbol{\beta}, \frac{1}{n} \mathbf{V})$.
2. $\mathbf{S}(\mathbf{x}) \in PSD(p)$, $\forall \mathbf{x} \in \mathbb{R}^{n \times p}$.
3. $\mathbf{S}(\mathbf{X}) = \mathbf{X}^T (\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T) \mathbf{X} \sim W_p(n-1, \mathbf{V})$
4. $\frac{1}{n-1} \mathbf{S}(\mathbf{X})$ ist ein erwartungstreuer Schätzer für \mathbf{V} ($n \geq 2$ vorausgesetzt).
5. Zusätzlich lässt sich zeigen, dass $\mathbf{S}, \bar{\mathbf{X}}$ unabhängig sind.

Theorem 1.7 Als ML-Schätzung ($\hat{\boldsymbol{\beta}}_{ML}(\mathbf{x}), \hat{\mathbf{V}}_{ML}(\mathbf{x})$) für $(\boldsymbol{\beta}, \mathbf{V})$ ergibt sich (wobei wir $n > p$ voraussetzen):

$$\hat{\boldsymbol{\beta}}_{ML}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \bar{\mathbf{x}}, \quad (1.17)$$

$$\hat{\mathbf{V}}_{ML}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \frac{1}{n} \mathbf{S}(\mathbf{x}) \quad (1.18)$$

Theorem 1.8 1) Als BLUE-Schätzung für $\boldsymbol{\beta}$ ergibt sich: $\hat{\boldsymbol{\beta}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \bar{\mathbf{x}}$.

2) Als MCUE-Schätzung (engl. minimum coriance unbiased estimator MCUE) für \mathbf{V} ergibt sich :

$$\hat{\mathbf{V}}(\mathbf{x}_1, \dots, \mathbf{x}_n) = \frac{1}{n-1} \mathbf{S}(\mathbf{x}). \quad (1.19)$$

(i) Der mittlere quadratische Fehler für $\hat{\boldsymbol{\beta}}$ ist gegeben durch:

$$\begin{aligned}MSE(\bar{\mathbf{X}}, \boldsymbol{\beta}) &= R(L(\bar{\mathbf{X}}, \boldsymbol{\beta})) \\ &= E_{\theta}[(\bar{\mathbf{X}} - \boldsymbol{\beta})^T (\bar{\mathbf{X}} - \boldsymbol{\beta})] \\ &= \text{Spur}(\text{Cov}(\bar{\mathbf{X}})) \\ &= \frac{1}{n} \text{Spur}(\mathbf{V})\end{aligned}$$

(ii) Der gewichtete quadratische Fehler für $\widehat{\boldsymbol{\beta}}$, mit der gegebenen Gewichtsfunktion

$$\mathbf{W} : \mathbb{R}^p \times PD(p) \rightarrow PD(p), \mathbf{W}(\boldsymbol{\beta}, \mathbf{V}) = \mathbf{V}^{-1} \quad (1.20)$$

ist gegeben durch:

$$\begin{aligned} MSE(\bar{\mathbf{X}}, \boldsymbol{\beta}) &= R(L(\bar{\mathbf{X}}, \boldsymbol{\beta})) \\ &= E_{\boldsymbol{\theta}}[(\bar{\mathbf{X}} - \boldsymbol{\beta})^T \mathbf{V}^{-1} (\bar{\mathbf{X}} - \boldsymbol{\beta})] \\ &= \text{Spur}(\mathbf{V}^{-1} \text{Cov}(\bar{\mathbf{X}})) \\ &= \frac{p}{n}. \end{aligned}$$

Modell mit multivariaten gemeinsam, u.i.v. normalverteilten ZV'en

Wir betrachten das Modell mit $n \geq (p + q)$ u.i.v. \mathbb{R}^{p+q} -wertigen, normal-verteilten Zufallsvariablen, wobei der $((p + q)$ -dim.) Erwartungswert und die (positiv definite $(p + q) \times (p + q)$) Kovarianzmatrix die Parameter sind:

$$\begin{aligned} \mathbf{X}_i &\sim N_p(\boldsymbol{\beta}_x, \mathbf{V}_{xx}), \quad \mathbf{Y}_i \sim N_q(\boldsymbol{\beta}_y, \mathbf{V}_{yy}) \\ \begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix} &\sim N_{p+q}(\boldsymbol{\beta}, \mathbf{V}), \quad (1 \leq i \leq n) \text{ u.i.v.}, \\ \boldsymbol{\beta} &= \begin{pmatrix} \boldsymbol{\beta}_x \\ \boldsymbol{\beta}_y \end{pmatrix} \in \mathbb{R}^{p+q} \text{ und } \mathbf{V} = \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \in PD(p+q). \end{aligned}$$

Wir haben also effektiv den Parameter $\boldsymbol{\theta} = (\boldsymbol{\beta}_x, \boldsymbol{\beta}_y, \mathbf{V}_{xx}, \mathbf{V}_{yy}, \mathbf{V}_{xy})$ mit fünf Komponenten und den Parameterbereich:

$$\Theta = \{ \boldsymbol{\theta} = (\boldsymbol{\beta}_x, \boldsymbol{\beta}_y, \mathbf{V}_{xx}, \mathbf{V}_{yy}, \mathbf{V}_{xy}) : \boldsymbol{\beta}_x \in \mathbb{R}^p, \boldsymbol{\beta}_y \in \mathbb{R}^q, \mathbf{V}_{xx} \in PD(p), \mathbf{V}_{yy} \in PD(q), \mathbf{V}_{xy} \in \mathbb{R}^{p \times q} \text{ mit } \mathbf{V}_{yy} - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{yx}^T \in PD(q) \}.$$

Wir betrachten die Statistiken:

$$\begin{aligned} \bar{\mathbf{x}} &= \frac{1}{n} \mathbf{x}^T \mathbf{1}_n, & \bar{\mathbf{y}} &= \frac{1}{n} \mathbf{y}^T \mathbf{1}_n & \text{mit } \mathbf{x} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^{n \times q} \\ \mathbf{S}_{xx} &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T, & \mathbf{S}_{yy} &= \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^T \\ \mathbf{S}_{xy} &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{y}_i - \bar{\mathbf{y}})^T, & \mathbf{S}_{yx} &= \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{x}_i - \bar{\mathbf{x}})^T. \end{aligned}$$

Mit Hilfe von Theorem 1.2 erhält man:

Theorem 1.9 1) Als BLUE-Schätzung $(\widehat{\boldsymbol{\beta}}_x, \widehat{\boldsymbol{\beta}}_y)$ für $(\boldsymbol{\beta}_x, \boldsymbol{\beta}_y)$ ergibt sich :

$$\widehat{\boldsymbol{\beta}}_x(\mathbf{x}) = \bar{\mathbf{x}} \quad \widehat{\boldsymbol{\beta}}_y(\mathbf{y}) = \bar{\mathbf{y}}.$$

2) Als MCUE-Schätzung $(\widehat{\mathbf{V}}_{xx}, \widehat{\mathbf{V}}_{yy}, \widehat{\mathbf{V}}_{xy})$ für $(\mathbf{V}_{xx}, \mathbf{V}_{yy}, \mathbf{V}_{xy})$ ergibt sich:

$$\begin{aligned} \widehat{\mathbf{V}}_{xx}(\mathbf{x}) &= \frac{1}{n-1} \mathbf{S}_{xx}, \quad \widehat{\mathbf{V}}_{yy}(\mathbf{y}) = \frac{1}{n-1} \mathbf{S}_{yy} \\ \widehat{\mathbf{V}}_{xy}(\mathbf{x}, \mathbf{y}) &= \frac{1}{n-1} \mathbf{S}_{xy}, \quad \text{mit } \mathbf{x} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^{n \times q}. \end{aligned}$$

1.6. Autoregressives Modell

In der Praxis spielen autoregressive Modelle eine wichtige Rolle, bei denen sich der aktuelle Messwert als Funktion des vorherigen Messwertes und eines Fehlerterms ergibt. Im Folgenden betrachten wir das autoregressive Modell von der Ordnung 1, kurz AR(1). Für Beweise und weitere Ergebnisse auf die Literatur, z.B. Kreiß und Neuhaus (2006, Kap.2, Kap.3)

Definition 1.9 (AR(1)-Modell) Die Komponenten der \mathbb{R}^p -wertigen Zufallsvariable \mathbf{Y} sind autoregressiv AR(1), falls:

$$(i) \quad Y_1 = e_1 \sim N_1(0, \frac{1}{1-\rho^2}\sigma^2) \quad |\rho| < 1.$$

(ii) $Y_i = \rho Y_{i-1} + e_i$, mit $e_i \sim N_1(0, \sigma^2)$, für $i = 2, \dots, p$ und e_1, \dots, e_p unabhängig sind.

Interpretation: Der neue Wert Y_i hängt direkt vom unmittelbaren Vorgänger Y_{i-1} ab,

$$\begin{aligned} Y_i &= e_i + \rho Y_{i-1} \\ &= e_i + \rho(\rho Y_{i-2} + e_{i-1}) \\ &= e_i + \rho e_{i-1} + \rho^2 Y_{i-2} \end{aligned}$$

Eigenschaften des AR(1)-Modells

1. **Erwartungswert:** $E(Y_i) = 0$.
2. **Varianz:** $Var(Y_i) = \frac{1}{1-\rho^2}\sigma^2$.
3. **Autokovarianz:** $Cov(Y_i, Y_{i-1}) = \rho Var(Y_{i-1}) = \rho \frac{\sigma^2}{1-\rho^2}$
und für $s \leq i - 1$ gilt auch:

$$Cov(Y_i, Y_{i-s}) = \rho^s \frac{\sigma^2}{1-\rho^2}.$$

Im Rahmen dieser Arbeit betrachten wir die standardisierte Version der Kovarianzmatrix mit $\sigma^2 = 1 - \rho^2$, z.B. für $p = 4$ erhalten wir:

$$Cov(\mathbf{Y}) = \begin{pmatrix} 1 & \rho & \rho^2 & \rho^3 \\ \rho & 1 & \rho & \rho^2 \\ \rho^2 & \rho & 1 & \rho \\ \rho^3 & \rho^2 & \rho & 1 \end{pmatrix}.$$

2. Assoziationsstrukturen

In diesem Abschnitt werden Regressionsmodelle und Korrelationskoeffizienten für ein-dimensionale Zufallsvariable Y und einen Zufallsvektor \mathbf{X} behandelt. Danach kann die Kanonische Korrelation als Verallgemeinerung der Korrelationskoeffizienten aufgefasst werden. Daran schließt sich der wichtige Spezialfall unter Annahme der Normalverteilung an.

2.1. Theorie der Regressionsmodelle

In dem folgenden Abschnitt wird kurz in grundlegende Prinzipien der Regressionsmodelle eingeführt. Speziell der multiple Korrelationskoeffizient wird dann definiert. Wir verweisen für Beweise und weitere Ergebnisse auf die Literatur, z.B. Puntanen, Styan und Isotalo (2011, Kap.9) und Rao (1973, Kap.4g) zur Regressionsmodelle und Fujikoshi, Ulyanov und Shimizu (2010, Kap.4.2) und Alvin (2008, Kap.10.3) zum multiplen Korrelationskoeffizient.

Modell: In diesem Abschnitt betrachten wir das Modell 1.3.a mit $q = 1$. Kurz:

$$\begin{aligned} & \begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix} \text{ eine } \mathbb{R}^{p+1} \text{ wertige ZV, mit} \\ & \boldsymbol{\mu} \in \mathbb{R}^{p+1}, \mathbf{V} \in PD(p+1) \\ & \boldsymbol{\mu}, \mathbf{V} \text{ bekannt} \end{aligned}$$

Notation:

- Bezeichne

$$SP = \{g | g : \mathbb{R}^p \rightarrow \mathbb{R}, g \text{ messbar und } E(g^2(\mathbf{X})) < \infty\}$$

die Menge aller Funktionen g , so dass $g(\mathbf{X})$ eine quadratintegrierbare reelle Zufallsvariable auf \mathbb{R}^p ist.

-

$$LP = \{g | g : \mathbb{R}^p \rightarrow \mathbb{R}, g(\mathbf{x}) = b + \mathbf{a}^T \mathbf{x}, \mathbf{a} \in \mathbb{R}^p, b \in \mathbb{R}\} \subseteq SP$$

die Menge aller linearen reellwertigen Transformationen ist.

Interessieren wir uns für den Zusammenhang zwischen der Zufallsvariable Y und Zufallsvektor \mathbf{X} . So gibt sich es verschiedene Ansätze, diesen zu modellieren. Einer der bekanntesten ist das Regressionsmodell.

2. Assoziationsstrukturen

Definition 2.1 Eine Funktion $m \in SP$ heißt Regressionsfunktion von Y unter \mathbf{X} , wenn gilt

$$E[(Y - m(\mathbf{X}))^2] = \min_{g \in SP} E[(Y - g(\mathbf{X}))^2].$$

Definition 2.2 Eine lineare Funktion $f \in LP$ heißt lineare Regressionsfunktion von Y unter \mathbf{X} , wenn gilt

$$E[(Y - f(\mathbf{X}))^2] = \min_{l \in LP} E[(Y - l(\mathbf{X}))^2].$$

Bestimmung der Regressionsfunktion

Aus der Definition des bedingten Erwartungswertes folgt

Satz 2.1 Rao (1973, Kap.4g.1) Sei $\hat{Y} = m(\mathbf{X}) = E(Y|\mathbf{X})$, dann ist $E(Y - g(\mathbf{X}))^2$ minimal, wenn $g(\mathbf{X}) = m(\mathbf{X})$ ist.

Nun fragen wir nach der besten linearen Vorhersage von Y unter \mathbf{X} in der linearen Klasse LP

Satz 2.2 Rao (1973, Kap.4g.1) Sei $\hat{Y} = f(\mathbf{X}) = \mu_y + \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}(\mathbf{X} - \boldsymbol{\mu}_x)$, dann ist $E(Y - l(\mathbf{X}))^2$ minimal, wenn $l(\mathbf{X}) = f(\mathbf{X})$ ist

und es gilt:

Satz 2.3 In der Klasse der linearen Regressionsmodelle gilt, dass

$$\text{Var}(Y - \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}\mathbf{X}) \leq \text{Var}(Y - \mathbf{a}^T\mathbf{X}), \quad \forall \mathbf{a} \in \mathbb{R}^p \quad (2.1)$$

ist.

Bemerkung 2.1 $\hat{\beta} = \mathbf{V}_{xx}^{-1}\mathbf{V}_{xy}$ heißt Regressionskoeffizient von Y unter \mathbf{X} .

Die Abweichungen zwischen den tatsächlichen Werte Y und den geschätzten Werten \hat{Y} :

$$\epsilon = Y - \hat{Y} \quad \text{bezeichnet man als Residuen.}$$

Statistische Eigenschaften der Residuen:

Wir beenden diesen Teil mit der Untersuchung der statistischen Eigenschaften der Residuen $\epsilon = Y - E(Y|\mathbf{X})$

1. Erwartungswert: $E(\epsilon) = E(Y - E(Y|\mathbf{X})) = 0$, d.h. die Residuen sind im Mittel Null.

2. Varianz: Es gilt

$$\begin{aligned}
 \text{Var}(\epsilon) &= \text{Var}(Y - E(Y|\mathbf{X})) \\
 &= \text{Var}(Y) + \text{Var}(E(Y|\mathbf{X})) - 2 \text{Cov}(Y, E(Y|\mathbf{X})) \\
 &= \text{Var}(Y) + \text{Var}(E(Y|\mathbf{X})) - 2 \text{Cov}(E(Y|\mathbf{X}), E(Y|\mathbf{X})) \\
 &= \text{Var}(Y) - \text{Var}(E(Y|\mathbf{X})).
 \end{aligned}$$

In der Klasse der linearen Regression erhalten wir:

$$\text{Var}(\epsilon) = \text{Var}(Y - f(\mathbf{X})) \quad (2.2)$$

$$= \text{Var}(Y - \mu_y - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1}(\mathbf{X} - \boldsymbol{\mu}_x)) \quad (2.3)$$

$$= V_{yy} - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \quad (2.4)$$

Hieraus erhalten wir:

Satz 2.4 Die folgenden Bedingungen sind äquivalent

1. $\text{Var}(Y - \mathbf{b}^T \mathbf{X}) \geq \text{Var}(Y - \mathbf{a}^T \mathbf{X}), \forall \mathbf{b} \in \mathbb{R}^p$
2. $\text{Cov}(Y - \mathbf{a}^T \mathbf{X}, \mathbf{c}^T \mathbf{X}) = 0, \forall \mathbf{c} \in \mathbb{R}^p.$
3. $\mathbf{a} = \mathbf{V}_{xx}^{-1} \mathbf{X}_{xy}.$

Im Kapitel 3 werden wir die Eigenschaften der Residuen der Vorhersage in verschiedenen Klassen untersuchen.

Multipler Korrelationskoeffizient

Eine wichtige Maßzahl zur Beurteilung der Güte einer Regressionsfunktion ist der multiple Korrelationskoeffizient

Definition 2.3 (Multipler Korrelationskoeffizient)

Der multiple Korrelationskoeffizient $\rho_{Y|\mathbf{X}}$ ist bestimmt als die Korrelation zwischen Y und Regressionskoeffizient von Y unter \mathbf{X} und gegeben durch:

$$\rho_{Y|\mathbf{X}} = \frac{\text{Cov}(\boldsymbol{\beta}^T \mathbf{X}, Y)}{\sqrt{\text{Var}(\boldsymbol{\beta}^T \mathbf{X})} \sqrt{\text{Var}(Y)}}$$

Eigenschaften des multiplen Korrelationskoeffizients

1. Es gilt:

$$\text{Var}(\mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X}) = \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} = \text{Cov}(\mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X}, Y) \quad (2.5)$$

Aus (2.3) und (2.4) ergibt sich die Zerlegung

$$\text{Var}(Y) = \text{Var}(\mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X}) + \text{Var}(\epsilon)$$

2. Assoziationsstrukturen

2. Der quadrierte multiple Korrelationskoeffizient stellt sich damit in der Form

$$\rho_{Y|\mathbf{X}}^2 = \frac{\text{Var}(\mathbf{V}_{y\mathbf{x}}\mathbf{V}_{\mathbf{xx}}^{-1}\mathbf{X})}{\text{Var}(Y)} = 1 - \frac{\text{Var}(\epsilon)}{\text{Var}(Y)}$$

dar

3. Der multiple Korrelationskoeffizient $\rho_{Y|\mathbf{X}}$ ist die maximale Korrelation zwischen Y und einer Linearkombination $\mathbf{a}^T\mathbf{X}$

$$\rho_{Y|\mathbf{X}} = \max_{\mathbf{a} \in \mathbb{R}^p}(\mathbf{a}^T\mathbf{X}, Y) \quad (2.6)$$

4. Aus (2.3) und (2.4) erhalten wir

$$\text{Var}(\epsilon) = V_{yy} - \mathbf{V}_{y\mathbf{x}}\mathbf{V}_{\mathbf{xx}}^{-1}\mathbf{V}_{\mathbf{x}y} \quad (2.7)$$

$$= \text{Var}(Y)(1 - \rho_{Y|\mathbf{X}}^2), \quad (2.8)$$

d.h. je größer $\rho_{Y|\mathbf{X}}^2$, desto kleiner ist $\text{Var}(\epsilon)$, das bedeutet in diesem Sinn, dass $\rho_{Y|\mathbf{X}}^2$ ein Maß für die Güte der Vorhersage von Y unter \mathbf{X} durch die Regressionsfunktion ist.

Insgesamt erhalten wir somit das folgende Ergebnis, dass der multiple Korrelationskoeffizient die Stärke des Zusammenhangs einer Zufallsvariable Y mit einem Zufallsvektor \mathbf{X} erfasst. In Kapitel 3 werden wir auf die Begriffe Vorhersage und multipler Korrelationskoeffizient zurückkommen.

2.2. Kanonische Korrelation

Im Folgenden geben wir einen kurzen Überblick über die wichtigsten Resultate zum Begriff der kanonischen Korrelation und verweisen für Beweise und weitere Ergebnisse auf die Literatur, z.B. Schlittgen (2009, Kap.12) und Fujikoshi, Ulyanov und Shimizu (2010, Kap.11). Der Ansatz der kanonischen Korrelation bietet eine Möglichkeit Zusammenhänge zwischen zwei Gruppen von Variablen herzuleiten. Hierbei bestimmen wir ein Paar Linearkombinationen der beiden Variablen, so dass diese die größte mögliche Korrelation aufweisen. Dies ist die Aufgabe der sogenannten kanonischen Korrelationsanalyse, die von Hotelling in den 30er Jahren entwickelt wurde.

Die Ausgangsfragestellung der kanonischen Korrelation

Im letzten Abschnitt wurde der Korrelationskoeffizienten für eine eindimensionale Variable Y und einen \mathbb{R}^p -Zufallsvektor \mathbf{X} dargestellt. Die kanonische Korrelation kann in folgendem Sinn als Verallgemeinerung aufgefasst werden: Nun ist auch \mathbf{Y} ein \mathbb{R}^q Zufallsvektor ($q > 1$). Die Idee der kanonischen Korrelationsanalyse ist die Analyse der Struktur des gemeinsamen Zusammenhangs von zwei Merkmalsgruppen in niedrigerer

Dimension mit Hilfe von Projektionen.

Wir betrachten wieder in diesem Abschnitt das Modell 1.3.a mit $q > 1$, d.h.

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \text{ eine } \mathbb{R}^{p+q}\text{-wertige ZV, mit} \\ \boldsymbol{\mu} \in \mathbb{R}^{p+q}, \mathbf{V} \in PD(p+q) \\ \boldsymbol{\mu}, \mathbf{V} \text{ bekannt.}$$

Gesucht sind Linearkombinationen $\mathbf{a}^T \mathbf{X}$ und $\mathbf{b}^T \mathbf{Y}$ mit $\mathbf{a} \in \mathbb{R}^p \setminus \{0\}$, $\mathbf{b} \in \mathbb{R}^q \setminus \{0\}$ derart, dass die Korrelation zwischen diesen Linearkombinationen maximal wird. Dazu betrachtet man

$$\text{Corr}(\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}) = \frac{\text{Cov}(\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y})}{\sqrt{\text{Var}(\mathbf{a}^T \mathbf{X})} \sqrt{\text{Var}(\mathbf{b}^T \mathbf{Y})}} = \frac{\mathbf{a}^T \mathbf{V}_{xy} \mathbf{b}}{\sqrt{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \sqrt{\mathbf{b}^T \mathbf{V}_{yy} \mathbf{b}}}$$

Vereinbarung:

- die (größte) kanonische Korrelation zwischen X_1, \dots, X_p und Y_1, \dots, Y_q ist die maximale Korrelation zwischen allen Linearkombinationen $\mathbf{a}^T \mathbf{X}$ und $\mathbf{b}^T \mathbf{Y}$,

$$\text{Corr}_1(\mathbf{X}, \mathbf{Y}) = \max_{\mathbf{a} \in \mathbb{R}^p \setminus \{0\}, \mathbf{b} \in \mathbb{R}^q \setminus \{0\}} \text{Corr}(\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}).$$

- die Linearkombinationen $\mathbf{U} = \mathbf{a}^T \mathbf{X}$ und $\mathbf{V} = \mathbf{b}^T \mathbf{Y}$, für welche das Maximum angenommen wird, heißen kanonische Variablen.

Das Ziel ist, die Gewichtungsvektoren \mathbf{a}, \mathbf{b} zu bestimmen unter den Normierungsbedingungen $\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a} = 1$, $\mathbf{b}^T \mathbf{V}_{yy} \mathbf{b} = 1$.

Bestimmung der Gewichtungsvektoren

Kleines Hilfsresultat Seien \mathbf{A} eine $q \times p$ -Matrix, \mathbf{B} eine $p \times q$ -Matrix. Die Matrizen \mathbf{AB} und \mathbf{BA} haben die gleichen von Null verschiedenen Eigenwerte.

Es lässt sich zeigen: die Gewichtungsvektoren $\mathbf{a}_{cc}, \mathbf{b}_{cc}$ ($cc = \text{canonical correlation}$), sind die Lösung von

$$(\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{V}_{yy}^{-1} \mathbf{V}_{yx} - r^2 \mathbf{I}_p) \mathbf{a} = \mathbf{0}_{p \times p} \quad \text{und} \quad (\mathbf{V}_{yy}^{-1} \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} - r^2 \mathbf{I}_q) \mathbf{b} = \mathbf{0}_{q \times q},$$

wobei r der größte Eigenwert von

$$\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{V}_{yy}^{-1} \mathbf{V}_{yx} \quad \text{bzw.} \quad \mathbf{V}_{yy}^{-1} \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \quad \text{ist.}$$

Beachte $\mathbf{A} = \mathbf{V}_{yy}^{-1} \mathbf{V}_{yx}$, $\mathbf{B} = \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy}$. Dann

$$\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{V}_{yy}^{-1} \mathbf{V}_{yx}, \quad \mathbf{V}_{yy}^{-1} \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy}$$

2. Assoziationsstrukturen

besitzen denselben Eigenwert.

Die Gewichtungsvektoren sind die zugehörigen Eigenvektoren, die zu dem größten Eigenwert r gehören.

$$\mathbf{a}_{cc} = e.v_{max}(\mathbf{V}_{xx}^{-1}\mathbf{V}_{xy}\mathbf{V}_{yy}^{-1}\mathbf{V}_{yx}) \quad \text{und} \quad \mathbf{b}_{cc} = e.v_{max}(\mathbf{V}_{yy}^{-1}\mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}\mathbf{V}_{xy})$$

und es gilt auch:

$$1) \quad \mathbf{a}_{cc} = \mathbf{V}_{xx}^{-1}\mathbf{V}_{xy}\mathbf{b}_{cc} \quad \text{und} \quad \mathbf{b}_{cc} = \mathbf{V}_{yy}^{-1}\mathbf{V}_{yx}\mathbf{a}_{cc}.$$

2) Das Minimum $\min_{\mathbf{a}, \mathbf{b}} E[(\mathbf{a}^T \mathbf{X} - \mathbf{b}^T \mathbf{Y})^2]$ wird angenommen für

$$\mathbf{a}_{cc} = \mathbf{V}_{xx}^{-1}\mathbf{V}_{xy}\mathbf{b}_{cc} \quad \text{und} \quad \mathbf{b}_{cc} = \mathbf{V}_{yy}^{-1}\mathbf{V}_{yx}\mathbf{a}_{cc}.$$

Andere Zugänge zur Bestimmung der Gewichtungsvektoren \mathbf{a}, \mathbf{b} :

Das Optimierungsproblem

$$\max_{\mathbf{a}, \mathbf{b}} \text{Corr}^2(\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y})$$

kann in zwei Schritten erfolgen:

1: wenn \mathbf{b} fest ist, erhalten wir

$$\begin{aligned} \max_{\mathbf{a}} \text{Corr}^2(\mathbf{a}^T \mathbf{X}, \mathbf{b}^T \mathbf{Y}) &= \max_{\mathbf{a}} \frac{(\mathbf{a}^T \mathbf{V}_{xy} \mathbf{b})^2}{(\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a})(\mathbf{b}^T \mathbf{V}_{yy} \mathbf{b})} \\ &= \frac{1}{\mathbf{b}^T \mathbf{V}_{yy} \mathbf{b}} \max_{\mathbf{a}} \frac{(\mathbf{a}^T \mathbf{V}_{xy} \mathbf{b})^2}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \\ &= \frac{\mathbf{b}^T \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{b}}{\mathbf{b}^T \mathbf{V}_{yy} \mathbf{b}} \quad (\text{nach Satz A.4}). \end{aligned}$$

2: danach

$$\max_{\mathbf{b}} \frac{\mathbf{b}^T \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{b}}{\mathbf{b}^T \mathbf{V}_{yy} \mathbf{b}} = \lambda_{max}(\mathbf{V}_{yy}^{-1} \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{yx}^T) \quad (\text{nach Satz A.3}).$$

Eigenschaften der kanonischen Korrelation

Die kanonische Korrelation hat folgende Eigenschaften:

1. Ihr Wertebereich ist bestimmt durch

$$0 \leq \text{Corr}_1(\mathbf{X}, \mathbf{Y}) \leq 1$$

2. Im Fall $q = 1$ und $p > 1$ ergibt sich der multiple Korrelationskoeffizient zwischen Y und $\mathbf{X} = (X_1, \dots, X_p)$

3. Für $\mathbf{a} \in \mathbb{R}^p \setminus \{0\}$ und $\mathbf{b} \in \mathbb{R}^q \setminus \{0\}$ gilt :

$$E[(\mathbf{a}^T \mathbf{X} - \mathbf{b}^T \mathbf{Y})^2] \geq 2(1 - \text{Corr}_1^2(\mathbf{Y}, \mathbf{X}))$$

unter den Nebenbedingungen: $\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a} = 1, \mathbf{b}^T \mathbf{V}_{yy} \mathbf{b} = 1.$

4. Wenn $\mathbf{Y}^* = \mathbf{F}\mathbf{Y} + \mathbf{f}$ und $\mathbf{X}^* = \mathbf{G}\mathbf{X} + \mathbf{g}$ sind, wobei $\mathbf{G} : p \times p$ und $\mathbf{F} : q \times q$ reguläre Matrizen und $\mathbf{f} : p \times 1$ und $\mathbf{g} : q \times 1$ feste Vektoren sind, dann ist die kanonische Korrelation zwischen \mathbf{Y}^* und \mathbf{X}^* die gleiche wie zwischen \mathbf{Y} und \mathbf{X} , das heißt die kanonische Korrelation ist invariant unter allen linearen Transformationen, die vollen Rang haben.

Bemerkung 2.2 Wenn $q = 1$ ist, dann ist der multiple Korrelationskoeffizient zwischen Y und \mathbf{X}^* der gleiche wie zwischen Y und \mathbf{X} , das heißt der multiple Korrelationskoeffizient ist invariant unter allen linearen Transformationen von \mathbf{X} die vollen Rang haben.

2.3. Normalverteilungsmodell

Betrachte des Modells 1.3.b. Kurz:

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim N_{p+q} \left(\begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix}, \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \right),$$

$$\begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix} \in \mathbb{R}^{p+q}, \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \in PD(p+q).$$

Aus den Verteilungsergebnissen für (mehrdimensionale) normalverteilte Zufallsvariablen und aus der Modelldefinition des Regressionsmodells ergibt sich

- Die Regressionsfunktion von \mathbf{Y} , die auf \mathbf{X} basiert, ist

$$m(\mathbf{X}) = E(\mathbf{Y}|\mathbf{X}) = \boldsymbol{\mu}_y + \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} (\mathbf{X} - \boldsymbol{\mu}_x)$$

- Die Regressionsfunktion schreibt sich als lineare Funktion um

$$\begin{aligned} E(\mathbf{Y}|\mathbf{X}) &= \boldsymbol{\mu}_y + \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} (\mathbf{X} - \boldsymbol{\mu}_x) \\ &= \boldsymbol{\mu}_y - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \boldsymbol{\mu}_x + \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X} \\ &= \boldsymbol{\alpha}_0 + \boldsymbol{\alpha}_1^T \mathbf{X}, \end{aligned}$$

wobei

$$\begin{aligned} \boldsymbol{\alpha}_0 &= \boldsymbol{\mu}_y - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \boldsymbol{\mu}_x \\ \boldsymbol{\alpha}_1 &= \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \end{aligned}$$

sind.

2. Assoziationsstrukturen

- Die Kovarianzmatrix von $(\mathbf{Y} - E(\mathbf{Y}|\mathbf{X}))$ ist minimal und gegeben durch:

$$\begin{aligned}\boldsymbol{\Sigma} &= \text{Cov}(\mathbf{Y} - E(\mathbf{Y}|\mathbf{X})) \\ &= E((E(\mathbf{Y}|\mathbf{X}) - \mathbf{Y})(E(\mathbf{Y}|\mathbf{X}) - \mathbf{Y})^T) \\ &= \mathbf{V}_{yy} - \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}\mathbf{V}_{xy}\end{aligned}$$

- Der mittlere quadratische Fehler der Regressionsfunktion von \mathbf{Y} , die auf \mathbf{X} basiert, ist gegeben durch den Ausdruck:

$$\begin{aligned}MSE(\mathbf{Y}, E(\mathbf{Y}|\mathbf{X})) &= E((\mathbf{Y} - E(\mathbf{Y}|\mathbf{X}))^T(\mathbf{Y} - E(\mathbf{Y}|\mathbf{X}))) \\ &= \text{Spur}(\text{Cov}(\mathbf{Y} - E(\mathbf{Y}|\mathbf{X}))) \\ &= \text{Spur}(\mathbf{V}_{yy} - \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}\mathbf{V}_{xy})\end{aligned}$$

- Verteilung der Residuen $\mathbf{Y} - m(\mathbf{X})$ mit Normalverteilungsannahme:

$$\boldsymbol{\epsilon} \sim N_q(\mathbf{0}, \mathbf{V}_{yy} - \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}\mathbf{V}_{xy}) \quad (2.9)$$

- Verteilung der Residuenquadratsumme:

$$\boldsymbol{\epsilon}^T \mathbf{D}^{-1} \boldsymbol{\epsilon} \sim \chi_q^2, \text{ mit } \mathbf{D} = \mathbf{V}_{yy} - \mathbf{V}_{yx}\mathbf{V}_{xx}^{-1}\mathbf{V}_{xy} \in PD(q). \quad (2.10)$$

Schätzung von $\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_1, \boldsymbol{\Sigma}$

Wir betrachten das Modell 1.5 (Modell mit multivariaten gemeinsam, u.i.v. normalverteilten ZV'en): Kurz

$$\begin{aligned}\begin{pmatrix} \mathbf{X}_i \\ \mathbf{Y}_i \end{pmatrix} &\sim N_{p+q}(\boldsymbol{\beta}, \mathbf{V}), \quad (1 \leq i \leq n) \quad u.i.v., \\ \boldsymbol{\beta} &= \begin{pmatrix} \boldsymbol{\beta}_x \\ \boldsymbol{\beta}_y \end{pmatrix} \in \mathbb{R}^{p+q} \text{ und } \mathbf{V} = \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \in PD(p+q).\end{aligned}$$

Hieraus und aus Theorem 1.10 ergibt sich als MCUE-Schätzung (engl. minimum covariance unbiased MCUE) für $\boldsymbol{\alpha}_0, \boldsymbol{\alpha}_1$ und $\boldsymbol{\Sigma}$ nach Rao, Shalabh, Toutenburg und Heumann (2008, Theorem 3.16 S.74-75)

$$\begin{aligned}\hat{\boldsymbol{\alpha}}_0(\mathbf{x}, \mathbf{y}) &= \bar{\mathbf{y}} - \mathbf{S}_{yx}\mathbf{S}_{xx}^{-1}\bar{\mathbf{x}} \\ \hat{\boldsymbol{\alpha}}_1(\mathbf{x}, \mathbf{y}) &= \mathbf{S}_{xx}^{-1}\mathbf{S}_{xy} \\ \hat{\boldsymbol{\Sigma}}(\mathbf{x}, \mathbf{y}) &= \frac{1}{n-1}(\mathbf{S}_{yy} - \mathbf{S}_{yx}\mathbf{S}_{xx}^{-1}\mathbf{S}_{xy}) \\ \widehat{MSE}(\mathbf{Y}, E(\mathbf{Y}|\mathbf{X})) &= \text{Spur}\left[\frac{1}{n-1}(\mathbf{S}_{yy} - \mathbf{S}_{yx}\mathbf{S}_{xx}^{-1}\mathbf{S}_{xy})\right],\end{aligned}$$

mit $\mathbf{x} \in \mathbb{R}^{n \times p}, \mathbf{y} \in \mathbb{R}^{n \times q}$.

3. Einige Aspekte der Theorie der Vorhersage

Eine der wichtigsten Anwendungen in der Statistik ist die Prognose. Das folgende Kapitel beschäftigt sich mit der Optimalität von Vorhersagen. Hierfür wird der klassische Zugang der Effizienz, welche am mittleren quadratischen Fehler von dem zur schätzenden Vorhersage gemessen wird, betrachtet. Es stellt sich heraus, dass zusätzlich zu einem Abstandskriterium eine zweite Bedingung, die Unverzerrtheit, gefordert werden muss, um hinreichend allgemeine Aussagen treffen zu können. Anschließend werden als Erweiterung die Eigenschaften des mittleren quadratischen Fehlers behandelt, welche auch für die Darstellung der Methode der partiellen kleinsten Quadrate anwendbar sind.

3.1. Allgemeine Resultate

Das Ziel dieses Abschnitts wird sein, ein Kriterium zu finden, welches bereits zur Beurteilung einer Vorhersage genutzt werden kann, um die Qualität einer Vorhersage zu messen.

Modell: Seien $(\Omega, \mathcal{C}, \mathbb{P})$ ein W-Raum und

$$\begin{aligned} X &: (\Omega, \mathcal{C}) \rightarrow (M_1, \mathcal{A}_1) \\ Y &: (\Omega, \mathcal{C}) \rightarrow (M_0, \mathcal{A}_0) \end{aligned}$$

zwei Zufallsvariablen gegeben, wobei

- (M_i, \mathcal{A}_i) , $i = 0, 1$ Messräume sind.
- Y nicht beobachtbar ist.
- Die Verteilung $P^{(Y,X)}$ wird als gemeinsame Verteilung der Zufallsvariablen (Y, X) bezeichnet, und ist ein W-Verteilung auf der Produkt Sigma-Algebra $\mathcal{A}_0 * \mathcal{A}_1$.

In der Statistik ist die Frage von Interesse: wie lässt sich Y durch X approximativ beschreiben? oder anders gesagt, im Fall, dass Y nicht beobachtbar ist, welche Prognosen lassen sich über Y machen, wenn X beobachtbar ist?

Definition 3.1 (Vorhersage) Rao, Shalabh, Toutenburg und Heumann (2008, Kap.6)
Unter einer Vorhersage $\hat{Y} = g(X)$ für Y verstehen wir eine Funktion $g : M_1 \rightarrow M_0$, die

3. Einige Aspekte der Theorie der Vorhersage

für jeden möglichen Wert $x \in M_1$ der Zufallsvariable X einen Vorhersagewert $g(x)$ von Y liefert, dabei soll g messbar sein:

$$g : (M_1, \mathcal{A}_1) \rightarrow (M_0, \mathcal{A}_0).$$

Bezeichne

$$SP = \{g | g : M_1 \rightarrow M_0 \text{ messbar, } E(g^2(X)) < \infty\} \quad (3.1)$$

die Menge aller quadratintegrierbaren Vorhersagen von Y unter X .

Mit der Bezeichnung \widehat{Y} statt g für eine Vorhersage für Y haben wir einen kleinen Notationswechsel vorgenommen, der natürlich inhaltlich Bedeutung hat. In diesem Kapitel betrachten wir verschiedene Vorhersagen für die selbe unbeobachtbare Zufallsvariable Y , und bezeichnen diese mit g, g_1, \widehat{Y}, Y^* etc. für verschiedene Vorhersagen für Y . Das Konzept zur Beurteilung bzw. zum Vergleich verschiedener Vorhersagen für Y geht von einer gegebenen quadratischen Verlustfunktion L aus (steht für engl. Loss):

$$L : M_0 \times M_0 \rightarrow [0, \infty) \text{ messbar .}$$

Bemerkung 3.1 *Im Rahmen dieser Arbeit betrachten wir die Verlustfunktionen*

$$L : M_0 \times M_0 \rightarrow [0, \infty),$$

die die folgenden Bedingungen erfüllen:

1. $L \geq 0$
2. die Funktion L ist konvex.
3. $L(y_1, y_1) = 0$.

Speziell wird als Verlustfunktion verwendet $L(y_1, y_2) = \phi(y_1 - y_2)$, mit einer konvexen und Nullpunkt-symmetrischen Funktion ϕ .

Einige Verlustfunktionen:

(i) Im Fall $M_1 = \mathbb{R}^p, M_0 = \mathbb{R}$:

$$L(y_1, y_2) = (y_1 - y_2)^2 \quad (\text{quadratische Abweichung})$$

$$L(y_1, y_2) = |y_1 - y_2|, \quad (\text{absolute Abweichung})$$

(ii) Im Fall $M_1 = \mathbb{R}^p, M_0 = \mathbb{R}^q$:

$$L(\mathbf{y}_1, \mathbf{y}_2) = (\mathbf{y}_1 - \mathbf{y}_2)^T (\mathbf{y}_1 - \mathbf{y}_2) \quad (\text{quadratische Abweichung})$$

$$L_{\mathbf{W}}(\mathbf{y}_1, \mathbf{y}_2) = (\mathbf{y}_1 - \mathbf{y}_2)^T \mathbf{W} (\mathbf{y}_1 - \mathbf{y}_2), \quad \mathbf{W} \in PD(q) (\text{gewichtete quadratische Abweichung}).$$

Für das theoretische Verständnis der Vorhersage bzw. der partiellen kleinsten Quadrate spielt der mittlere quadratische Fehler der Vorhersage eine wichtige Rolle. Dazu führen wir folgende Definition ein.

Definition 3.2 (Risikofunktion einer Vorhersage) Rao, Shalabh, Toutenburg und Heumann (2008, Kap.6.4):

Für eine Vorhersage $\hat{Y} = g(X) : (M_1, \mathcal{A}_1) \rightarrow (M_0, \mathcal{A}_0)$ heißt

$$R(Y, \hat{Y}(X)) = E(L(Y, \hat{Y}(X))) \quad (3.2)$$

Risikofunktion der Vorhersage \hat{Y} .

Definition 3.3 (mittlerer quadratischer Fehler der Vorhersage MSEP)

(i) Für die quadratische Verlustfunktion $L(Y, \hat{Y}) = (Y - \hat{Y})^2$ heißt die Risikofunktion

$$R(Y, \hat{Y}) = E((Y - \hat{Y})^2), \quad (3.3)$$

der mittlere quadratische Fehler der Vorhersage: $MSEP(Y, \hat{Y})$, (engl. Mean Squared Error of Prediction)

(ii) Für die gewichtete quadratische Verlustfunktion

$$L_{\mathbf{W}}(\mathbf{Y}, \hat{\mathbf{Y}}) = (\mathbf{Y} - \hat{\mathbf{Y}})^T \mathbf{W} (\mathbf{Y} - \hat{\mathbf{Y}}), \quad \mathbf{W} \in PD(p) \quad \text{ist:}$$

$$R_{\mathbf{W}}(\mathbf{Y}, \hat{\mathbf{Y}}) = E((\mathbf{Y} - \hat{\mathbf{Y}})^T \mathbf{W} (\mathbf{Y} - \hat{\mathbf{Y}})) \quad (\text{Weighted Mean Squared Error}). \quad (3.4)$$

der gewichtete mittlere quadratische Fehler der Vorhersage (engl. Weighted Mean Squared Error)

Der mittlere quadratische Fehler der Vorhersage (MSEP) kann als Vergleichskriterium für verschiedene Vorhersagen herangezogen werden. Vorhersagen mit kleinem MSEP sind dabei vorzuziehen.

Folgerung 3.1 Für den mittleren quadratischen Fehler der Vorhersage gilt:

$$\begin{aligned} MSEP(Y, \hat{Y}) &= E[(Y - \hat{Y})^2] \\ &= E[(Y - \hat{Y} - E(Y - \hat{Y}) + E(Y - \hat{Y}))^2] \\ &= E[(Y - \hat{Y} - E(Y - \hat{Y}))^2] + (E(Y - \hat{Y}))^2 \\ &\quad + 2 E[E(Y - \hat{Y})(Y - \hat{Y})] - 2 E(Y - \hat{Y})E(Y - \hat{Y}) \\ &= E(Y - \hat{Y} - E(Y - \hat{Y}))^2 + (E(Y - \hat{Y}))^2 \\ &= \text{Var}(Y - \hat{Y}) + [E(Y) - E(\hat{Y})]^2. \end{aligned}$$

3. Einige Aspekte der Theorie der Vorhersage

Daraus erhalten wir folgende wichtige Zerlegung des mittleren quadratischen Fehlers der Vorhersage in die Varianz der Vorhersage und das Quadrat der Verzerrung:

$$\begin{aligned} R(Y, \hat{Y}) &= MSEP(Y, \hat{Y}) \\ &= \text{Var}(Y - \hat{Y}) + [E(Y) - E(\hat{Y})]^2. \end{aligned}$$

Im mehrdimensionalen Fall formulieren wir den mittleren quadratischen Fehler der Vorhersage wie folgt:

$$\begin{aligned} R(\mathbf{Y}, \hat{\mathbf{Y}}) &= MSEP(\mathbf{Y}, \hat{\mathbf{Y}}) \\ &= E((\mathbf{Y} - \hat{\mathbf{Y}})^T (\mathbf{Y} - \hat{\mathbf{Y}})) \\ &= \text{Spur} [\text{Cov}(\mathbf{Y} - \hat{\mathbf{Y}})] + [E(\mathbf{Y} - \hat{\mathbf{Y}})]^T [E(\mathbf{Y} - \hat{\mathbf{Y}})]. \end{aligned}$$

Definition 3.4 (optimale Vorhersage) Rao, Shalabh, Toutenburg und Heumann (2008, Kap.6.4)

Die Vorhersage $g^* \in G$ heißt optimale Vorhersage in der Klasse $G \subseteq SP$, wenn für jede Vorhersage $g \in G$ gilt:

$$\begin{aligned} R(Y, g^*(X)) &\leq R(Y, g(X)) \text{ bzw.} \\ R_{\mathbf{W}}(\mathbf{Y}, g^*(\mathbf{X})) &\leq R_{\mathbf{W}}(\mathbf{Y}, g(\mathbf{X})), \quad \mathbf{W} \in PD(p). \end{aligned}$$

Definition 3.5 (erwartungstreue Vorhersage)

die Vorhersage $g(X)$ heißt (engl. unbiased) erwartungstreue (unverzerrte) Vorhersage, wenn $E(g(X)) = E(Y)$ ist.

3.2. Optimalitätskonzepte

Das Ziel dieses Abschnittes besteht darin, die beste Vorhersage bzw. beste lineare Vorhersage zu bestimmen. Zu Beginn werden die folgenden Lemmas vorgestellt, welche die nachfolgenden Begriffe motivieren.

Lemma 3.1 (Minimaleigenschaft des Mittelwerts)

Sei X quadratintegrierbare reelle Zufallsvariable auf M_1 mit $\text{Var}(X) < \infty$ und $E(X) = \mu$. Dann gilt

$$E[(X - a^*)^2] = \min_{a \in \mathbb{R}} E[(X - a)^2] \Leftrightarrow a^* = \mu,$$

die Varianz ist ein Maß dafür, wie weit die Werte von X im Schnitt auseinander fallen

Lemma 3.2 Sei X quadratintegrierbare reelle Zufallsvariable auf M_1 mit $\text{Var}(X) < \infty$ und $E(X) = \mu$. Dann gilt

$$\left[\begin{array}{c} \text{für } a \in \mathbb{R} \text{ und} \\ E[(X - a)^2] = E(X^2) - a^2 \end{array} \right] \Leftrightarrow \left[\begin{array}{c} a = 0, \\ \text{oder} \\ a = \mu \end{array} \right].$$

Beweis:

" \Leftarrow " wenn $a = \mu$, folgt $E(X - \mu)^2 = E(X^2) - \mu^2$.

" \Rightarrow " Sei $E[(X - a)^2] = E(X^2) - a^2$,

$$E[(\mu - a)^2] = E(X^2) + a^2 - 2a\mu,$$

Hieraus und aus $E[(X - a)^2] = E(X^2) - a^2$ ergibt sich:

$$a(a - \mu) = 0 \implies a = \mu \quad \text{oder} \quad a = 0$$

und die Varianz ist gerade die mittlere quadratische Abweichung vom Mittelwert. Um Optimalitätsaussagen machen zu können, braucht man die folgenden Begriffe:

Definition 3.6 (Beste Vorhersage) Rao (1973, Kap. 4g.1)

Die Vorhersage $\hat{Y} \in G$ heißt beste Vorhersage von Y unter \mathbf{X} in der Klasse $G \subseteq SP$ im Sinne des mittleren quadratischen Fehlers, wenn gilt:

$$E(Y - \hat{Y}(X))^2 = \min_{g \in G} E(Y - g(\mathbf{X}))^2.$$

Einige Klassen von G :

1. $SP = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, g \text{ messbar } E(g^2(\mathbf{X})) < \infty\}$, die Menge aller quadratintegrierbaren reellen Zufallsvariablen
2. $UP = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, E(Y) - E(g(\mathbf{X})) = 0, E(g^2(\mathbf{X})) < \infty\} \subseteq SP$, die Teilmenge der unverzerrten quadratintegrierbaren reellen Zufallsvariablen
3. $UP_0 = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, E(g(\mathbf{X})) = 0, E(g^2(\mathbf{X})) < \infty\} \subseteq SP$, die Menge der quadratintegrierbaren reellen Zufallsvariablen mit Erwartungswert null.
4. $G_0 = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, g(\mathbf{x}) = b, b \in \mathbb{R}\} \subseteq SP$, die Menge aller konstanten Abbildungen
5. $LP = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, g(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b, \mathbf{a} \in \mathbb{R}^p, b \in \mathbb{R}\} \subseteq SP$, die Menge aller linearen reellwertigen Transformationen.
6. $LUP = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, g(\mathbf{x}) = \mathbf{a}^T \mathbf{x} + b, E(Y) - E(g(\mathbf{X})) = 0, \mathbf{a} \in \mathbb{R}^p, b \in \mathbb{R}\} \subseteq LP$ die Menge aller unverzerrten linearen reellwertigen Transformationen
7. $LP_2 = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, g(\mathbf{x}) = \mu_y + \mathbf{a}^T (\mathbf{x} - \boldsymbol{\mu}_x), \mathbf{a} \in \mathbb{R}^p, \text{ mit } \mathbf{a}^T \mathbf{V}_{xx} \mathbf{a} = 1\} \subseteq LP$ die Menge aller unverzerrten, über $Cov(\mathbf{X})$ standardisierten, linearen, reellwertigen Transformationen.
8. $LP_1 = \{g|g : \mathbb{R}^p \rightarrow \mathbb{R}, g(\mathbf{x}) = \mu_y + \mathbf{a}^T (\mathbf{x} - \boldsymbol{\mu}_x), \mathbf{a} \in \mathbb{R}^p \text{ mit } \mathbf{a}^T \mathbf{a} = 1\} \subseteq LP$, die Menge aller unverzerrten, standardisierten, linearen, reellwertigen Transformationen.

3. Einige Aspekte der Theorie der Vorhersage

Beachte: Folgende Aspekte sollte man beachten:

- 1) $LUP = LP \cap UP$
- 2) $G \cap UP$ die Teilmenge von G der unverzerrten quadratintegrierbaren reellen Zufallsvariablen
- 3) $G \cap UP_0$ die Teilmenge der quadratintegrierbaren reellen Zufallsvariablen mit Erwartungswert null.

Bezeichnungen: Eine gewisse Abschwächung der Behauptungen von Definition (3.6) werden wir im Folgenden erklären:

- (1) Eine Vorhersage $\hat{Y} \in SP$ mit $E[(Y - \hat{Y})^2] = \inf\{E[(Y - g(\mathbf{X}))^2] : g \in SP\}$, heißt beste Vorhersage.
- (2) Eine Vorhersage $\hat{Y} \in UP$ mit $Var(Y - \hat{Y}) = \inf\{Var(Y - g(\mathbf{X})) : g \in UP\}$, heißt beste erwartungstreue (unverzerrte) Vorhersage.
- (3) Eine Vorhersage $\hat{Y} \in G \cap UP$ mit $Var(Y - \hat{Y}) = \inf\{Var(Y - g(\mathbf{X})) : g \in G \cap UP\}$, heißt beste erwartungstreue (unverzerrte) Vorhersage in der Klasse G .
- (4) Eine Vorhersage $\hat{Y} \in LUP$ mit $Var(Y - \hat{Y}) = \inf\{Var(Y - g(\mathbf{X})) : g \in LUP\}$, heißt beste lineare erwartungstreue (unverzerrte) Vorhersage.
- (5) Seien G_1, G_2 zwei Klassen von G , wenn $G_1 \subset G_2$ ist, gilt:

$$\begin{aligned} \min_{h \in G_2} E(Y - h(\mathbf{X}))^2 &\leq \min_{g \in G_1} E(Y - g(\mathbf{X}))^2, \\ \max_{h \in G_2} Corr^2(Y, h(\mathbf{X})) &\geq \max_{g \in G_1} Corr^2(Y, g(\mathbf{X})) \end{aligned}$$

3.3. Die beste erwartungstreue Vorhersage

Betrachtet man nur erwartungstreue (unverzerrte) Vorhersagen, so kann man die Varianz der Vorhersage als Maß für die Qualität der Vorhersage heranziehen, da unter Erwartungstreue die Varianz der Vorhersage gleich dem mittleren quadratischen Fehler ist. Eine Vorhersage ist in diesem Sinn besser als alle anderen unverzerrten Vorhersagen, falls seine Varianz minimal auf G ist. Für eine erwartungstreue (unverzerrte) Vorhersage $\hat{Y} \in G \cap UP$ von Y gilt die folgende Optimalitätseigenschaft:

$$E[(Y - \hat{Y}(X))^2] = \min_{g \in G} E[(Y - g(\mathbf{X}))^2], \text{ d.h.} \quad (3.5)$$

$$Var(Y - \hat{Y}(X)) = \min_{g \in G \cap UP} Var(Y - g(\mathbf{X})). \quad (3.6)$$

BUP steht für beste erwartungstreue Vorhersage (engl. Best Unbiased Prediction).

Die Eigenschaften der besten erwartungstreuen Vorhersage

Die Optimalität von erwartungstreuen Vorhersagen lässt sich mit der Kovarianzmethode charakterisieren.

Lemma 3.3 (Kovarianzmethode)

Sei \hat{Y} eine erwartungstreue Vorhersage für Y in der Klasse G ($\hat{Y} \in G \cap UP$), dann gilt:

$$\hat{Y} \text{ ist BUP in } G \iff Cov(Y - \hat{Y}(\mathbf{X}), g(\mathbf{X})) = 0, \forall g \in G \cap UP_0.$$

Auf Grund der Modellannahme gilt folgende erweiterte Version der Kovarianzmethode

Theorem 3.1 (Kovarianzmethode-Version 2)

Sei \hat{Y} eine erwartungstreue Vorhersage für Y in der Klasse G ($\hat{Y} \in G \cap UP$), dann gilt:

$$\hat{Y} \text{ ist BUP in } G \iff Cov(Y - \hat{Y}(\mathbf{X}), g(\mathbf{X})) = 0, \forall g \in G.$$

Beweis: Notwendig: " \Rightarrow "

Seien \hat{Y} eine BUP für Y in der Klasse G und $g \in G$, dann gilt $\forall \lambda \in \mathbb{R}$

$$\begin{aligned} MSEP(Y, \hat{Y} - \lambda g(\mathbf{X})) &= Var(Y - \hat{Y} + \lambda g(\mathbf{X})) + \lambda^2 E^2(g(\mathbf{X})) \\ &= Var(Y - \hat{Y}) + \lambda^2 Var(g(\mathbf{X})) \\ &\quad + 2\lambda Cov(Y - \hat{Y}, g(\mathbf{X})) + \lambda^2 E^2(g(\mathbf{X})) \\ &\geq Var(Y - \hat{Y}) \quad (\text{da } \hat{Y} \text{ BUP ist}) \end{aligned}$$

$$\begin{aligned} \Rightarrow \lambda^2 Var(g(\mathbf{X})) + 2 \lambda Cov(Y - \hat{Y}, g(\mathbf{X})) + \lambda^2 E(g^2(\mathbf{X})) &\geq 0 \\ \Rightarrow \lambda^2 (Var(g(\mathbf{X})) + E^2(g(\mathbf{X}))) + 2 \lambda Cov(Y - \hat{Y}, g(\mathbf{X})) &\geq 0 \\ \Rightarrow \lambda^2 E(g^2(\mathbf{X})) + 2 \lambda Cov(Y - \hat{Y}, g(\mathbf{X})) &\geq 0 \\ \Rightarrow Cov(Y - \hat{Y}, g(\mathbf{X})) = 0, \text{ da} & \end{aligned}$$

wenn $Cov(Y - \hat{Y}, g(\mathbf{X})) < 0$ ist, $\exists \lambda \in [0, -2 \frac{Cov(Y - \hat{Y}, g(\mathbf{X}))}{E(g^2(\mathbf{X}))}]$, so dass

$$\lambda^2 E(g^2(\mathbf{X})) + 2 \lambda Cov(Y - \hat{Y}, g(\mathbf{X})) < 0.$$

ist, und wenn $Cov(Y - \hat{Y}, g(\mathbf{X})) > 0$ ist, $\exists \lambda \in [-2 \frac{Cov(Y - \hat{Y}, g(\mathbf{X}))}{E(g^2(\mathbf{X}))}, 0]$, so dass

$$\lambda^2 E(g^2(\mathbf{X})) + 2 \lambda Cov(Y - \hat{Y}, g(\mathbf{X})) < 0.$$

ist.

Hieraus folgt, dass \hat{Y} nicht BUP ist, wenn $Cov(Y - \hat{Y}, g(\mathbf{X})) \neq 0$ ist, im Widerspruch zur Annahme, dass \hat{Y} BUP für Y in der Klasse G .

3. Einige Aspekte der Theorie der Vorhersage

Hinreichend: " \Leftarrow "

Sei nun umgekehrt $Cov(Y - \hat{Y}(\mathbf{X}), g(\mathbf{X})) = 0, \forall g \in G$, dann für $\tilde{Y} \in G \cap UP$ gilt:

$$Cov(Y - \hat{Y}, \hat{Y}) = 0 \quad Cov(Y - \hat{Y}, \tilde{Y}) = 0$$

Hieraus ergibt sich, dass

$$Cov(Y - \hat{Y}, \hat{Y} - \tilde{Y}) = 0$$

und

$$\begin{aligned} Var(Y - \tilde{Y}) &= Var(Y - \hat{Y} + \hat{Y} - \tilde{Y}) \\ &= Var(Y - \hat{Y}) + Var(\hat{Y} - \tilde{Y}) \\ &\quad + 2 Cov(Y - \hat{Y}, \hat{Y} - \tilde{Y}) \\ &= Var(Y - \hat{Y}) + Var(\hat{Y} - \tilde{Y}). \end{aligned}$$

Dieses impliziert, dass

$$Var(Y - \hat{Y}) \leq Var(Y - \tilde{Y}).$$

Also ist \hat{Y} beste lineare erwartungstreue Vorhersage (BUP) und damit folgt die Behauptung.

Beispiel 3.1

(a) Für $\hat{Y} \in UP$ gilt das Äquivalent:

$$Cov(Y - \hat{Y}, g(\mathbf{X})) = 0, \forall g \in SP \iff \hat{Y} = m(\mathbf{X}),$$

wobei $m(\mathbf{X}) = E(Y|X)$ ist.

" \Leftarrow " Wenn $\hat{Y} = m(\mathbf{X})$ ist, erhalten wir aus den Eigenschaften der bedingten Erwartung,

$$\begin{aligned} Cov(Y - m(\mathbf{X}), g(\mathbf{X})) &= Cov(Y, g(\mathbf{X})) - Cov(m(\mathbf{X}), g(\mathbf{X})), \forall g \in SP \\ &= Cov(m(\mathbf{X}), g(\mathbf{X})) - Cov(m(\mathbf{X}), g(\mathbf{X})) \\ &= 0, \end{aligned}$$

da

$$\begin{aligned} Cov(Y, g(\mathbf{X})) &= E[(Y - E(Y))(g(\mathbf{X}) - E(g(\mathbf{X})))] \\ &= E\left[(g(\mathbf{X}) - E(g(\mathbf{X})))E[(Y - E(Y))|\mathbf{X}]\right] \\ &= Cov(m(\mathbf{X}), g(\mathbf{X})) \end{aligned}$$

" \Rightarrow " Sei nun umgekehrt $\hat{Y} \in UP$ mit $Cov(Y - \hat{Y}, g(\mathbf{X})) = 0, \forall g \in SP$, dann ist $\hat{Y} = m(\mathbf{X})$. Also genügt zu zeigen, dass

$$E(m(\mathbf{X}) - \hat{Y}) = 0 \quad \text{und} \quad Var(m(\mathbf{X}) - \hat{Y}) = 0$$

ist.

$$\begin{aligned}
 \text{Var}(m(\mathbf{X}) - \hat{Y}) &= \text{Cov}(m(\mathbf{X}) - \hat{Y}, m(\mathbf{X}) - \hat{Y}) \\
 &= \text{Cov}(Y - \hat{Y} - (Y - m(\mathbf{X})), m(\mathbf{X}) - \hat{Y}) \\
 &= \text{Cov}(Y - \hat{Y}, m(\mathbf{X}) - \hat{Y}) - \text{Cov}(Y - m(\mathbf{X}), m(\mathbf{X}) - \hat{Y}) \\
 &= 0 - 0
 \end{aligned}$$

Dieses impliziert, dass $\hat{Y} = m(\mathbf{X})$ und das Theorem(3.1) in der SP erfüllt ist.

(b) Auf analoge Weise gilt für $\hat{Y} \in LP \cap UP$, dass:

$$\text{Cov}(Y - \hat{Y}, g(\mathbf{X})) = 0, \forall g \in LP \setminus G_0 \iff \hat{Y} = \mu_y + \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} (\mathbf{X} - \boldsymbol{\mu}_x).$$

Zum Beispiel:

- 1) $\text{Cov}(Y - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X}, \mathbf{X}) = \mathbf{0}$
- 2) $\text{Cov}(Y - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X}, \mathbf{V}_{xx}^{-1} \mathbf{X}) = \mathbf{0}$
- 3) $\text{Cov}(Y - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X}, X_1) = 0, \quad X_1 = \mathbf{e}_1^T \mathbf{X}.$

(c) Im Fall, dass Y und \mathbf{X} unabhängig sind, gilt:

$$\text{Cov}(Y - \hat{Y}, g(\mathbf{X})) = 0, \forall g \in LP \setminus G_0 \iff \hat{Y} = \mu_y.$$

Beispiel 3.2 AR(1)-Modell

Wir betrachten das Modell 1.6 (Normalverteilung mit autoregressiver Struktur AR(1))

$$\begin{pmatrix} X_1 \\ X_2 \\ Y \end{pmatrix} \sim N_3 \left(\begin{pmatrix} \mu_{x_1} \\ \mu_{x_2} \\ \mu_y \end{pmatrix}, \begin{pmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{pmatrix} \right) \quad \text{mit } |\rho| < 1.$$

Seien

$$\begin{aligned}
 G_1 &= \{g|g : \mathbb{R}^2 \rightarrow \mathbb{R}, g(x_1, x_2) = a_1 x_1 : a_1 \in \mathbb{R}\} \\
 G_2 &= \{g|g : \mathbb{R}^2 \rightarrow \mathbb{R}, g(x_1, x_2) = c_1 x_1 + c_2 x_2 : c_1, c_2 \in \mathbb{R}\}
 \end{aligned}$$

Die beste lineare erwartungstreue Vorhersage für Y unter \mathbf{X} in den Klassen G_1, G_2 und der zugehörige mittlere quadratische Fehler sind:

G	BLUP	MSEP
G_1	$\rho^2 X_1$	$1 - \rho^4$
G_2	ρX_2	$1 - \rho^2$

3. Einige Aspekte der Theorie der Vorhersage

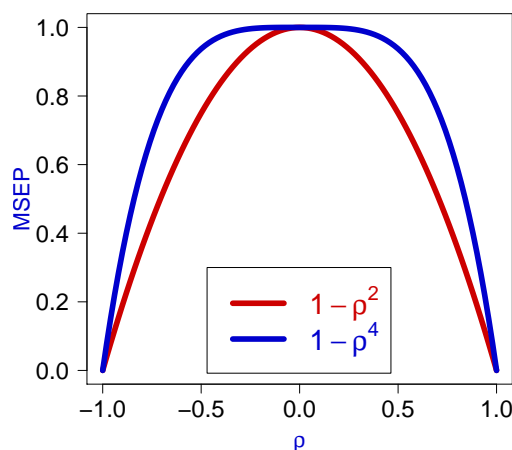


Abbildung 3.1.: Vergleich der MSEP

Die Abbildung (3.1) zeigt den mittleren quadratischen Fehler der Vorhersagen MSEP in den Klassen G_1 und G_2 in Abhängigkeit von ρ . Es gilt auch:

$$\begin{aligned} \text{Cov}(Y - \rho^2 X_1, aX_1) &= \text{Cov}(Y, aX_1) - \text{Cov}(\rho^2 X_1, aX_1), \forall a \in \mathbb{R} \\ &= \rho^2 a - \rho^2 a \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Cov}(Y - \rho X_2, c_1 X_1 + c_2 X_2) &= \text{Cov}(Y, c_1 X_1 + c_2 X_2) - \text{Cov}(\rho X_2, c_1 X_1 + c_2 X_2), \forall c_1, c_2 \in \mathbb{R} \\ &= \rho^2 c_1 + \rho c_2 - \rho^2 c_1 - \rho c_2 \\ &= 0. \end{aligned}$$

Weiterhin kann man bemerken:

$$\begin{aligned} \text{Cov}(Y - \rho^2 X_1, X_2) &= \text{Cov}(Y, X_2) - \text{Cov}(\rho^2 X_1, X_2) \\ &= \rho^2 - \rho^3 \end{aligned}$$

Beachte:

- Der Vorhersagefehler ist somit mit der für die Vorhersage zur Verfügung stehenden Information unkorreliert. In diesem Sinn kann die Vorhersage nicht mehr verbessert werden. Der Prognosefehler $Y - \hat{Y}$ steht somit orthogonal zu allen in die Prognosefunktion eingehenden Variablen X_1, \dots, X_p .
- Geometrisch ausgedrückt besagt diese Eigenschaft, dass man die beste lineare Vorhersage eines Punktes, hier y , in einem Teilraum, hier der durch x_1, \dots, x_p aufgespannten linearen Teilraum, als orthogonale Projektion des Punktes auf diesen Teilraum erhält.

- da $LP \subseteq SP$ ist, ist die lineare Vorhersage nicht immer optimal. Als Beispiel dafür ist:

Beispiel 3.3 Kreiß und Neuhaus (2006, Aufgabe(4.4))

Seien X, Z u.i.N(0, 1) -verteilt. Dann gilt für $Y = X^2 + Z$, dass

$$E(Y|X) = E((X^2 + Z)|X) = E(X^2|X) + E(Z) = X^2.$$

$$\hat{Y} = g(X) = BLUP(Y, X) = E(Y) = 1, \text{ da}$$

$$E(X) = 0, E(X^3) = 0 \text{ und } Cov(X^2, X) = 0.$$

Außerdem gilt für die Vorhersagefehler:

$$E(Y - E(Y|X))^2 = Var(Y - E(Y|X)) = 1.$$

$$E(Y - \hat{Y})^2 = E(Y - E(Y)) = Var(Y) = 3.$$

Auf ähnliche Weise gilt für $Y = X^2 + X + Z$, dass

$$E(Y|X) = E(X^2 + X + Z|X) = X^2 + X$$

$$\hat{Y}_1 = BLUP(Y, X) = E(Y) + \frac{Cov(Y, X)}{Var(X)}X = X + 1$$

$$E(Y - E(Y|X))^2 = 1 \quad E(Y - \hat{Y}_1)^2 = 3$$

sind.

Außerdem sind die folgenden Eigenschaften der Vorhersage nützlich, die sich unmittelbar aus dem Theorem (3.1) ergeben.

Folgerung 3.2 Wenn \hat{Y} eine erwartungstreue Vorhersage für Y in der Klasse G ist, dann gilt:

$$\hat{Y} \text{ ist BUP in } G \quad \iff \quad E[(Y - \hat{Y})g(\mathbf{X})] = 0, \quad \forall g \in G$$

Folgerung 3.3 Seien \tilde{Y}, \hat{Y} zwei beste erwartungstreue Vorhersagen für Y in der Klasse G , d.h. $(\tilde{Y}, \hat{Y} \in G \cap UP)$. Dann gilt: $\tilde{Y} = \hat{Y}$ fast sicher.

Beweis:

Kleines Hilfsresultat: $(a - b)^2 = a(a - b) - b(a - b)$

$$\begin{aligned} E[(\tilde{Y} - \hat{Y})^2] &= E[(Y - \hat{Y} - (Y - \tilde{Y}))^2] \\ &= E[(Y - \hat{Y})(\tilde{Y} - \hat{Y})] - E[(Y - \tilde{Y})(\tilde{Y} - \hat{Y})] \\ &= 0 \quad (\text{Folgerung 3.2}), \end{aligned}$$

hieraus ergibt sich, dass fast sicher $\tilde{Y} = \hat{Y}$ ist.

3. Einige Aspekte der Theorie der Vorhersage

Folgerung 3.4 Sei \hat{Y} BUP für Y in der Klasse $G \subseteq SP$, dann gilt:

$$\begin{aligned} \text{Var}(Y - \hat{Y}) &= \text{Cov}(Y - \hat{Y}, Y - \hat{Y}) \\ &= \text{Cov}(Y - \hat{Y}, Y) - \text{Cov}(Y - \hat{Y}, \hat{Y}) \\ &= \text{Cov}(Y - \hat{Y}, Y) \end{aligned}$$

Theorem 3.2 Sei \hat{Y} BUP für Y in der Klasse $G \subseteq SP$, dann gilt:

$$\text{Corr}^2(Y, \hat{Y}) = \frac{\text{Var}(\hat{Y})}{\text{Var}(Y)}.$$

Beweis:

Da \hat{Y} BUP ist, dann folgt $\text{Cov}(Y - \hat{Y}, \hat{Y}) = 0$, und

$$\begin{aligned} \text{Corr}^2(Y, \hat{Y}) &= \frac{\text{Cov}^2(\hat{Y}, Y)}{\text{Var}(\hat{Y})\text{Var}(Y)} \\ &= \frac{\text{Var}^2(\hat{Y})}{\text{Var}(\hat{Y})\text{Var}(Y)} \\ &= \frac{\text{Var}(\hat{Y})}{\text{Var}(Y)} \end{aligned}$$

Beispiel 3.4 Für $m(\mathbf{X}) = E(Y|\mathbf{X})$ gilt, dass

$$\min_{g \in SP} \text{Var}(Y - g(\mathbf{X})) = \text{Var}(Y - m(\mathbf{X})).$$

und aus den Eigenschaften der bedingten Erwartung erhalten wir

$$\begin{aligned} \text{Corr}^2(Y, m(\mathbf{X})) &= \frac{\text{Cov}^2(m(\mathbf{X}), Y)}{\text{Var}(m(\mathbf{X}))\text{Var}(Y)} \\ &= \frac{\text{Var}^2(m(\mathbf{X}))}{\text{Var}(m(\mathbf{X}))\text{Var}(Y)} \\ &= \frac{\text{Var}(m(\mathbf{X}))}{\text{Var}(Y)} \end{aligned}$$

Theorem 3.3 Sei \hat{Y} BUP in der Klasse $G \subseteq SP$, dann gilt

$$\text{Corr}^2(Y, \hat{Y}) = \max_{g \in G} \text{Corr}^2(Y, g(\mathbf{X})).$$

Beweis: Es gilt:

$$\begin{aligned}
 \text{Corr}^2(Y, g(\mathbf{X})) &= \frac{\text{Cov}^2(Y, g(\mathbf{X}))}{\text{Var}(Y)\text{Var}(g(\mathbf{X}))} \\
 &= \frac{\text{Cov}^2(\hat{Y}, g(\mathbf{X}))}{\text{Var}(\hat{Y})\text{Var}(g(\mathbf{X}))} \frac{\text{Var}(\hat{Y})}{\text{Var}(Y)} \quad (\text{Theorem 3.1}) \\
 &= \text{Corr}^2(\hat{Y}, g(\mathbf{X})) \frac{\text{Cov}^2(Y, \hat{Y})}{\text{Var}(\hat{Y})\text{Var}(Y)} \\
 &= \text{Corr}^2(\hat{Y}, g(\mathbf{X})) \text{Corr}^2(Y, \hat{Y}).
 \end{aligned}$$

Da $\text{Corr}^2(\hat{Y}, g(\mathbf{X})) \in [0, 1]$ ist, folgt

$$\text{Corr}^2(Y, \hat{Y}) \geq \text{Corr}^2(Y, g(\mathbf{X})).$$

3.4. Eigenschaften des mittleren quadratischen Fehlers

In diesem Abschnitt leiten wir verschiedene Formen für den mittleren quadratischen Fehler der Vorhersage her, welche für die nachfolgenden Begriffsbildungen und Aufgabenstellungen nützlich sind.

Folgerung 3.5 Sei \hat{Y} BUP in der Klasse $G \subseteq SP$, dann gilt

$$\text{MSEP}(Y, \hat{Y}) = \text{Var}(Y) - \text{Var}(\hat{Y})$$

Beweis:

$$\begin{aligned}
 \text{MSEP}(Y, \hat{Y}) &= \text{Var}(Y - \hat{Y}) \\
 &= \text{Var}(Y) + \text{Var}(\hat{Y}) - 2\text{Cov}(Y, \hat{Y}) \\
 &= \text{Var}(Y) + \text{Var}(\hat{Y}) - 2\text{Var}(\hat{Y}) \\
 &= \text{Var}(Y) - \text{Var}(\hat{Y})
 \end{aligned}$$

Folgerung 3.6 Sei \hat{Y} BUP für Y in der Klasse $G \subseteq SP$, dann

$$\text{MSEP}(Y, \hat{Y}) = \text{Var}(Y) - \text{Cov}(Y, \hat{Y}).$$

Folgerung 3.7 Sei \hat{Y} BUP für Y in der Klasse $G \subseteq SP$, dann

$$\text{MSEP}(Y, \hat{Y}) = \text{Var}(Y)(1 - \text{Corr}^2(Y, \hat{Y}))$$

3. Einige Aspekte der Theorie der Vorhersage

Beweis:

$$\begin{aligned} MSEP(Y, \hat{Y}) &= Var(Y - \hat{Y}) \\ &= Var(Y) - Corr^2(Y, \hat{Y})Var(Y) \quad (\text{Theorem 3.2}) \\ &= Var(Y)(1 - Corr^2(Y, \hat{Y})). \end{aligned}$$

Für die beste erwartungstreue Vorhersagen \hat{Y} in der Klasse $G \subseteq SP$ lässt sich anschaulich die Korrelation $Corr^2(Y, \hat{Y})$ interpretieren: $Corr^2(Y, \hat{Y})$ ist ein Maß dafür, um wie viel besser Y durch eine Zufallsvariable \hat{Y} angenähert werden kann (im Sinne des mittleren quadratischen Fehlers) als durch eine Vorhersage $g_0 \in G_0$. Es gilt nämlich

$$\min_{g \in G} E[(Y - g(X))^2] = (1 - Corr^2(Y, \hat{Y})) \min_{g_0 \in G_0} E[(Y - g_0(X))^2]$$

Zum Beweis: $E[(Y - g_0(X))^2] = Var(Y) + (E(Y) - g_0(X))^2$, hat $Var(Y)$ als minimalen Wert, während

$$\begin{aligned} E[(Y - g(X))^2] &= E[(Y - \hat{Y} + \hat{Y} - g(X))^2] \\ &= Var(Y)(1 - Corr^2(Y, \hat{Y})) + (g(X) - \hat{Y})^2 \quad (\text{Folgerung 3.7}) \end{aligned}$$

als Minimum $Var(Y)(1 - Corr^2(Y, \hat{Y}))$ für $\hat{Y} = g(X)$ annimmt.

Bemerkung 3.2 Aus Theorem 3.2 und Folgerung 3.7 ergibt sich, dass

$$Corr^2(Y, \hat{Y}) = 1 - \frac{Var(Y - \hat{Y})}{Var(Y)} \quad (3.7)$$

Interpretation:

- (i) $Corr^2(Y, \hat{Y})$ ist ein Maß der Genauigkeit der Vorhersage, das man beim Vergleich verschiedener Situationen oder verschiedener Klassen von Vorhersagen verwendet.
- (ii) Je näher der Ausdruck (3.7) bei 1 liegt, desto kleiner ist der Vorhersagefehler

$$Var(Y - \hat{Y})$$

d.h. desto besser ist die Vorhersage.

(iii) Im Extremfall gilt:

$$\begin{aligned} Corr^2(Y, \hat{Y}) = 1 &\iff Var(Y - \hat{Y}) = 0 \quad \text{d.h. } Y = \hat{Y} \text{ f.s. und} \\ &\text{die Vorhersage ist perfekt} \\ Corr^2(Y, \hat{Y}) = 0 &\iff Var(Y) = Var(Y - \hat{Y}) \text{ d.h. } \hat{Y} \in G_0 \text{ und} \\ &\text{die Vorhersage ist gleich dem Mittelwert } \hat{Y} = \mu_y \text{ f.s.} \end{aligned}$$

Beachte:

Wenn \hat{Y} BUP für Y auf G ist, ergibt sich hieraus und aus den Darstellungsformeln für den mittleren quadratischen Fehler, dass

a)

$$\text{Var}(Y) - \text{Var}(\hat{Y}) = \min_{g \in G} \text{MSEP}(Y, g(\mathbf{X})).$$

b)

$$\text{Var}(Y) - \text{Cov}(Y, \hat{Y}) = \min_{g \in G} \text{MSEP}(Y, g(\mathbf{X})).$$

c)

$$\text{Var}(Y)(1 - \text{Corr}^2(Y, \hat{Y})) = \min_{g \in G} \text{MSEP}(Y, g(\mathbf{X})).$$

Mit Hilfe der obigen Formeln lässt sich die optimale Lösung der Minimierung des mittleren quadrierten Vorhersagefehlers wie folgt schreiben:

a)

$$\text{Var}(Y) - \text{Var}(\hat{Y}) = \min_{g \in G \cap UP} \text{Var}(Y - g(\mathbf{X}))$$

b)

$$\text{Var}(Y) - \text{Cov}(Y, \hat{Y}) = \min_{g \in G \cap UP} \text{Var}(Y - g(\mathbf{X}))$$

c)

$$\text{Var}(Y)(1 - \text{Corr}^2(Y, \hat{Y})) = \min_{g \in G \cap UP} \text{Var}(Y - g(\mathbf{X})).$$

Eine nützliche Erweiterung obiger Begriffe ist die folgende Äquivalenz

Lemma 3.4 *Es gilt:*

$$\left[\begin{array}{l} \hat{Y}(\mathbf{X}) \text{ erwartungstreu in der Klasse } G \subseteq SP \text{ mit} \\ \text{MSEP}(Y, \hat{Y}(\mathbf{X})) = \text{Var}(Y) - \text{Var}(\hat{Y}(\mathbf{X})) \end{array} \right] \Leftrightarrow \left[\begin{array}{l} \hat{Y}(\mathbf{X}) = \mu_y \\ \text{oder} \\ \hat{Y}(\mathbf{X}) \text{ BUP} \end{array} \right].$$

3.5. Modell der Dimensionsreduktion

Das Ziel einer Dimensionsreduktion ist, eine Abbildung zu finden, die von einem höherdimensionalen Raum in einen niederdimensionalen Raum abbildet. So werden häufig Projektionen der Daten auf einen niedrigdimensionalen Unterraum betrachtet, wobei dann zu entscheiden ist, welche der möglichen Projektionen, z.B. alle zweidimensionalen Koordinaten, auszuwählen sind. Aus mathematischer Sicht hat eine Dimensionsreduktion häufig das Ziel, eine niedrigere Dimensionalität zu erreichen, indem eine Abbildung $f : \mathbb{R}^p \rightarrow \mathbb{R}^d, d < p$, benutzt wird. Dabei werden wir in dieser Arbeit nur endlich

3. Einige Aspekte der Theorie der Vorhersage

dimensionale Vektorräume betrachtet. Bei statistischen Fragestellungen kann unter Dimensionsreduktion sowohl eine Reduzierung der Anzahl von Variablen als auch eine Verringerung der Anzahl von Beobachtungen verstanden werden. Hier bezieht sich der Begriff Dimensionsreduktion auf die Anzahl von Variablen. Natürlich sollte eine Reduktion so durchgeführt werden, dass keine wesentlichen Informationen verloren gehen. Um zu einer mathematischen Präzisierung dieser Forderung zu gelangen, stellen wir zunächst folgende Überlegung an:

Modelldefinition

Wir betrachten in diesem Abschnitt weiter das Modell 1.3.a mit $q = 1$, $p > 1$, d.h.

$$\begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix} \text{ eine } \mathbb{R}^{p+1} \text{ wertige ZV, wobei} \\ \boldsymbol{\mu} \in \mathbb{R}^{p+1}, \mathbf{V} \in PD(p+1) \text{ und bekannt sind.}$$

Ein solches Modell für Dimensionsreduktion wird von Li (1991) und Cook (1994) verwendet und im Folgenden vorgestellt:

Zwischen einer Zielgröße Y und einem Vektor \mathbf{X} gelte der Zusammenhang:

$$\begin{aligned} Y &= d(Z_1, \dots, Z_k, \epsilon), \quad \text{mit } k < p, \\ &= d(\mathbf{Z}, \epsilon), \end{aligned}$$

so dass $Z_1 = \mathbf{a}_1^T \mathbf{X}, \dots, Z_k = \mathbf{a}_k^T \mathbf{X}$ sind. Die Matrixformulierung ist:

$$\mathbf{Z} = \mathbf{U}^T \mathbf{X} \quad \text{mit} \quad \mathbf{U} = [\mathbf{a}_1, \dots, \mathbf{a}_k],$$

wobei die folgenden Voraussetzungen vereinbart werden:

- (i) ϵ sei eine \mathbb{R} -wertige ZV, so dass \mathbf{X}, ϵ unabhängig seien.
- (ii) Die Kovarianz V_ϵ der Residuen ϵ sei unbekannt.
- (iii) Die Vektoren $\mathbf{a}_1, \dots, \mathbf{a}_k \in \mathbb{R}^p$ heißen dimensionsreduzierende Richtung (Abkürzung DR-Richtung), mit $\mathbf{U} = [\mathbf{a}_1, \dots, \mathbf{a}_k]$ eine $p \times k$ -Matrix.
- (iv) Die unbekannte Abbildung $d : \mathbb{R}^k \rightarrow \mathbb{R}$ werde als Linkfunktion der Reduktion bezeichnet.

Interpretation:

1. Der Spaltenraum von \mathbf{U} wird mit $\mathcal{U} = \text{span}(\mathbf{U})$ bezeichnet. Dieser Unterraum heißt dimensionsreduzierender Unterraum (DRU). Weiterhin sei angenommen, dass \mathcal{U} ein eindeutig bestimmter DRU minimaler Dimension k sei mit $1 \leq k \leq p$ d.h. man hofft, dass k klein ist.

2. Die Vektoren $\mathbf{a}_1, \dots, \mathbf{a}_k \in \mathbb{R}^p$ bzw. der von ihnen aufgespannte Raum $\mathcal{U} = \text{span}(\mathbf{U})$ sowie die minimale Anzahl k sind unbekannt.
3. Die Abbildung $d : \mathbb{R}^k \rightarrow \mathbb{R}$ ist unbekannt. Das primäre Ziel besteht also zunächst darin, eine Dimensionsreduktion zu erreichen, die die relevanten Informationen in \mathbf{X} repräsentiert. Aus mathematischer Sicht besteht das Ziel darin, einen Unterraum zu suchen, der durch $\mathbf{U} = [\mathbf{a}_1, \dots, \mathbf{a}_k]$ aufgespannt wird, so dass $k \ll p$. Da $\mathbf{U} = [\mathbf{a}_1, \dots, \mathbf{a}_k]$ nicht eindeutig bestimmt ist, wird vorausgesetzt, dass $\mathbf{a}_1, \dots, \mathbf{a}_k$ unabhängig sind und eine Orthonormalbasis von $\mathcal{U} = \text{span}(\mathbf{U})$ bilden, d.h. es gilt, dass $\mathbf{U}^T \mathbf{U} = \mathbf{I}_k$ ist.
4. Die Effizienz dieser Reduktion ist gegeben durch:

$$\text{Eff}(d(Z)) = \frac{\text{MSEP}(Y, \text{BLUP}(Y, \mathbf{X}))}{\text{MSEP}(Y, \text{BLUP}(Y, d(Z)))}$$

Beispiel 3.5 (AR(1)-Modell)

Wir betrachten das Modell 1.6 (Normalverteilung mit autoregressiver Struktur AR(1):)

$$\begin{pmatrix} X_1 \\ X_2 \\ Y \end{pmatrix} \sim N_3 \left(\begin{pmatrix} \mu_{x_1} \\ \mu_{x_2} \\ \mu_y \end{pmatrix}, \begin{pmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{pmatrix} \right) \quad \text{mit } |\rho| < 1.$$

Die Unterräume:

$$\mathcal{U}_1 = \text{span}((1, 0)^T), \quad \mathcal{U}_2 = \text{span}((0, 1)^T), \quad \mathcal{U}_3 = \text{span}\left(\left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}\right)^T\right)$$

sind dimensionsreduzierende Unterräume DRU.

- Das reduzierte Modell unter Unterraum \mathcal{U}_1 ist:

$$\begin{pmatrix} Z_1 \\ Y \end{pmatrix} \sim N_2 \left(\begin{pmatrix} \mu_{x_1} \\ \mu_y \end{pmatrix}, \begin{pmatrix} 1 & \rho^2 \\ \rho^2 & 1 \end{pmatrix} \right), \quad \text{mit}$$

$$Z_1 = X_1 \quad \text{und} \quad d(Z_1) = \text{BLUP}(Y, Z_1).$$

$$\text{MSEP}(Y, \text{BLUP}(Y, Z_1)) = 1 - \rho^4 \quad \text{und} \quad \text{Eff}(d(Z_1)) = \frac{1}{1 + \rho^2}.$$

- Das reduzierte Modell unter Unterraum \mathcal{U}_2 ist:

$$\begin{pmatrix} Z_2 \\ Y \end{pmatrix} \sim N_2 \left(\begin{pmatrix} \mu_{x_2} \\ \mu_y \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix} \right), \quad \text{mit}$$

$$Z_2 = X_2 \quad \text{und} \quad d(Z_2) = \text{BLUP}(Y, Z_2).$$

$$\text{MSEP}(Y, \text{BLUP}(Y, Z_2)) = 1 - \rho^2 \quad \text{und} \quad \text{Eff}(d(Z_1)) = 1$$

$$\text{Eff}(d(Z_1)) = 1 \quad (\text{Eigenschaften des AR(1)-Modells})$$

3. Einige Aspekte der Theorie der Vorhersage

- Das reduzierte Modell unter Unterraum \mathcal{U}_3 ist:

$$\begin{pmatrix} Z_3 \\ Y \end{pmatrix} \sim N_2 \left(\begin{pmatrix} \frac{1}{\sqrt{2}}(\mu_{x_1} + \mu_{x_2}) \\ \mu_y \end{pmatrix}, \begin{pmatrix} 1 + \rho & \frac{1}{\sqrt{2}}(\rho^2 + \rho) \\ \frac{1}{\sqrt{2}}(\rho^2 + \rho) & 1 \end{pmatrix} \right), \text{ mit}$$

$$Z_3 = \frac{1}{\sqrt{2}}(X_1 + X_2) \quad \text{und} \quad d(Z_3) = BLUP(Y, Z_3)$$

$$MSEP(Y, BLUP(Y, Z_3)) = 1 - \frac{\rho^2(1 + \rho)}{2}, \quad \text{Eff}(d(Z_3)) = \frac{1 - \rho^2}{1 - \frac{\rho^2(1 + \rho)}{2}}$$

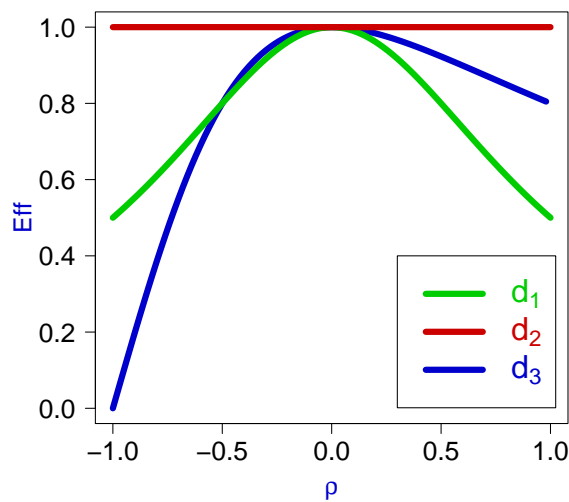


Abbildung 3.2.: Effizienz der Reduktionen nach d_1, d_2, d_3

Die Abbildung (3.2) zeigt die Effizienzfunktionen der Dimensionsreduktionen d_1, d_2 und d_3 in Abhängigkeit von ρ . Man sieht, falls ρ an der Stelle nahe bei 0 liegt, ist die Effizienz von d_1, d_3 nahe 1 und es gibt keinen Informationsverlust bei $d_2, \forall |\rho| < 1$.

4. Vorhersage in Linearen Modellen

Eines der wichtigsten Ziele der Regressionsrechnung ist es, den gefundenen Zusammenhang auszunutzen, um zukünftige Werte zu prognostizieren, wenn man eine neue Beobachtung für den Regressor erhält. Dabei geht man davon aus, dass der durch die Regressionsgerade beschriebene Zusammenhang auch in der Zukunft zumindest für einen gewissen Zeitraum Gültigkeit besitzt.

Wir beginnen mit dem folgenden Beispiel, um einige Ideen einzubringen und die nachfolgende Begriffsbildungen zu motivieren.

Beispiel 4.1 (*Stochastisch unabhängige identische reelle Zufallsvariablen*)

$$X_1, \dots, X_n, X_{n+1} \text{ u.i.v. } N(\beta, \sigma^2), \quad (4.1)$$

wobei, der Parameter $\theta = (\beta, \sigma^2)^T \in \mathbb{R} \times (0, \infty)$ unbekannt ist. Die Zufallsvariablen X_1, \dots, X_n seien beobachtbar und X_{n+1} nicht beobachtbar. Das Ziel ist, die Vorhersage von X_{n+1} zu bestimmen.

Maßtheoretische Modellformulierung:

$$(\mathbb{R}^{n+1}, \mathcal{B}^{n+1}, (P_{\beta, \sigma^2})_{(\beta, \sigma^2) \in \mathbb{R} \times (0, \infty)}), \text{ wobei } P_{\beta, \sigma^2} = [N(\beta, \sigma^2)]^n \otimes N(\beta, \sigma^2).$$

- Die beste lineare Vorhersage von X_{n+1} unter X_1, \dots, X_n ist

$$\begin{aligned} \hat{X}_{n+1} &= \hat{\beta} \\ &= \frac{1}{n} \mathbf{1}_n^T \mathbf{X} \\ &= \bar{X} \end{aligned}$$

- Damit gilt insbesondere:

$$\begin{aligned} \text{Var}(X_{n+1} - \hat{X}_{n+1}) &= \text{Var}(X_{n+1}) + \text{Var}(\hat{X}_{n+1}) - 2 \text{Cov}(X_{n+1}, \hat{X}_{n+1}) \\ &= \sigma^2 + \frac{\sigma^2}{n}. \end{aligned}$$

Nach diesem einführenden Beispiel kehren wir zur Modelldarstellung zurück.

4.1. Modelldefinition

Definition 4.1 (*Lineares Modell*)

Seien $p, n \in \mathbb{N}$ mit $p < n$. Ein lineares Modell für \mathbb{R}^n -wertige beobachtete Zufallsvariable \mathbf{X} mit unbekanntem $(p+1)$ -dimensionalen Parameter $\theta = (\beta^T, \sigma^2)^T \in \mathbb{R}^p \times (0, \infty)$ besteht aus:

4. Vorhersage in Linearen Modellen

- einer reellen bekannten $n \times p$ -Designmatrix \mathbf{F} von vollem Rang p , und
- einem Zufallsvektor $\mathbf{e}_x = (e_1, \dots, e_n)^T$ von n standardisierten Zufallsvariablen, den Fehlern.

Der n -dimensionale Beobachtungsvektor \mathbf{X} ergibt sich durch die Matrixformulierung:

$$\mathbf{X} = \mathbf{F}\beta + \mathbf{e}_x, \quad \text{mit} \quad \text{Cov}(\mathbf{e}_x) = \sigma^2 \mathbf{V}_{xx}, \quad \mathbf{V}_{xx} \text{ bekannt} \quad (4.2)$$

Das zugehörige statistische Experiment unter Normalverteilungsannahme ist:

$$(\mathbb{R}^n, \mathcal{B}^n, N_n(F\beta, \sigma^2 V)_{\theta \in \Theta}) : \Theta = \mathbb{R}^p \times (0, \infty), \quad \theta = (\beta^T, \sigma^2)^T.$$

Basierend auf der Definition des linearen Modells ergibt sich die Vorhersage im linearen Modell durch die Matrixformulierung: (siehe Puntanen, Styan und Isotalo (2011, Kap.10.9))

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{F} \\ \mathbf{F}_0 \end{pmatrix} \beta + \begin{pmatrix} \mathbf{e}_x \\ \mathbf{e}_y \end{pmatrix}, \quad (4.3)$$

mit

- \mathbf{F}_0 einer reellen bekannten $q \times p$ - Designmatrix von vollem Rang $q : q < p$
- \mathbf{Y} nicht beobachtbare (q -dimensionale) Zufallsvariable.

$$Y : (\Omega, \mathcal{C}) \rightarrow (\mathbb{R}^q, \mathcal{B}^q),$$

- $\mathbf{e}_x, \mathbf{e}_y$ gemeinsam verteilt, mit Kovarianzmatrix:

$$\Sigma = \text{Cov} \begin{pmatrix} \mathbf{e}_x \\ \mathbf{e}_y \end{pmatrix} = \sigma^2 \mathbf{V}, \quad \mathbf{V} = \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{xy}^T & \mathbf{V}_{yy} \end{pmatrix} \in PD(n+q) \text{ bekannt.}$$

Maßtheoretische Modellformulierung zum Beispiel unter Normalverteilungsannahme:

$$\left(\mathbb{R}^{n+q}, \mathcal{B}^{n+q}, \left(N_{n+q} \left(\begin{bmatrix} \mathbf{F}\beta \\ \mathbf{F}_0\beta \end{bmatrix}, \sigma^2 \begin{bmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{xy}^T & \mathbf{V}_{yy} \end{bmatrix} \right) \right)_{(\beta, \sigma^2) \in \mathbb{R}^p \times (0, \infty)} \right). \quad (4.4)$$

Definition 4.2 (lineare Vorhersage) Puntanen, Styan und Isotalo (2011, Kap.10.9)
Eine spezielle Form der linearen Vorhersage für \mathbf{Y} in diesem Modell ist:

$$\mathbf{Y}(\hat{\beta}) = \mathbf{F}_0 \hat{\beta} + \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} (\mathbf{X} - \mathbf{F} \hat{\beta}), \quad (4.5)$$

wobei $\hat{\beta}$ ein Schätzer für β ist, (engl. linear prediction (LP)).

Als Folgerung der Definition 4.2 ergibt sich:

Folgerung 4.1 (BLUP-Vorhersage) *Puntanen, Styan und Isotalo (2011, Kap.10.9) und Isotalo, Puntanen (2009)*

- 1) Eine lineare Vorhersage $\mathbf{Y}(\hat{\boldsymbol{\beta}})$ in der Form (4.5) heißt erwartungstreue Vorhersage, wenn $\hat{\boldsymbol{\beta}}$ erwartungstreuer Schätzer für $\boldsymbol{\beta}$ ist, (engl. linear unbiased prediction (LUP)).
- 2) Wenn $\hat{\boldsymbol{\beta}}$ ein BLUE-Schätzer ist, heißt

$$\mathbf{Y}(\hat{\boldsymbol{\beta}}) = \mathbf{F}_0 \hat{\boldsymbol{\beta}} + \mathbf{V}_{\mathbf{yx}} \mathbf{V}_{\mathbf{xx}}^{-1} (\mathbf{X} - \mathbf{F} \hat{\boldsymbol{\beta}}) \quad (4.6)$$

$$= BLUE(\mathbf{F}_0 \boldsymbol{\beta}) + \mathbf{V}_{\mathbf{yx}} \mathbf{V}_{\mathbf{xx}}^{-1} (\mathbf{X} - BLUE(\mathbf{F} \boldsymbol{\beta})) \quad (4.7)$$

die beste lineare erwartungstreue Vorhersage (engl. best linear unbiased prediction (BLUP)).

Einige Eigenschaften

Lemma 4.1 *Schmidt(1988)*

Seien $\mathbf{Z} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$, $\mathbf{G} = \begin{pmatrix} \mathbf{F} \boldsymbol{\beta} \\ \mathbf{F}_0 \boldsymbol{\beta} \end{pmatrix}$ und $\mathbf{V} = \begin{pmatrix} \mathbf{V}_{\mathbf{xx}} & \mathbf{V}_{\mathbf{xy}} \\ \mathbf{V}_{\mathbf{xy}}^T & \mathbf{V}_{\mathbf{yy}} \end{pmatrix}$, dann gilt:

$$E((\mathbf{Z} - \mathbf{G} \hat{\boldsymbol{\beta}})^T \mathbf{V}^{-1} (\mathbf{Z} - \mathbf{G} \hat{\boldsymbol{\beta}})) = E((\mathbf{X} - \mathbf{F} \hat{\boldsymbol{\beta}})^T \mathbf{V}_{\mathbf{xx}}^{-1} (\mathbf{X} - \mathbf{F} \hat{\boldsymbol{\beta}})) +$$

$$E((\mathbf{Y} - \mathbf{F}_0 \hat{\boldsymbol{\beta}})^T (\mathbf{V}_{\mathbf{yy}} - \mathbf{V}_{\mathbf{yx}} \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}})^{-1} (\mathbf{Y} - \mathbf{F}_0 \hat{\boldsymbol{\beta}})),$$

wenn $\hat{\boldsymbol{\beta}}$ ein linearer erwartungstreuer Schätzer ist.

Lemma 4.2 (Kovarianz der Residuen) *Sengupta und Jammalamadaka (2003, Kap.7.13.1)*

Wenn $\hat{\boldsymbol{\beta}}$ ein linearer erwartungstreuer Schätzer ist, dann hat der mittlere quadratische Fehler der Vorhersage die folgende Form

$$\begin{aligned} Cov(\mathbf{Y} - \mathbf{Y}(\hat{\boldsymbol{\beta}})) &= \sigma^2 (\mathbf{V}_{\mathbf{yy}} - \mathbf{V}_{\mathbf{yx}} \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}}) \\ &\quad + (\mathbf{F}_0 - \mathbf{V}_{\mathbf{yx}} \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{F}) Cov(\hat{\boldsymbol{\beta}}) (\mathbf{F}_0^T - \mathbf{F}^T \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}}) \end{aligned}$$

Im Hinblick auf die Definition 4.2 gilt die folgende spezielle Version der Kovarianzmethode (siehe Lemma 3.3).

Lemma 4.3 (Kovarianzmethode)

Seien $\tilde{\boldsymbol{\beta}}$, $\boldsymbol{\beta}^*$ und $\boldsymbol{\beta}^{**}$ lineare erwartungstreue Schätzer, dann gilt:

$$Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) = \mathbf{0}_{q \times q} \Leftrightarrow \mathbf{Y}(\tilde{\boldsymbol{\beta}}) \text{ ist BLUP,}$$

4. Vorhersage in Linearen Modellen

Beweis: (Für eindimensionale Aspekte siehe Beweis von Theorem 3.1)

Notwendig: Seien $\mathbf{Y}(\tilde{\boldsymbol{\beta}})$ BLUP, dann ist $\mathbf{Y}(\tilde{\boldsymbol{\beta}}) - \lambda(\mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**}))$ LUP für \mathbf{Y} , $\forall \lambda \in \mathbb{R}$ mit

$$E(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}) + \lambda(\mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**}))) = \mathbf{0}_{q \times 1}.$$

Hieraus ergibt sich, dass

$$\begin{aligned} Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}) + \lambda(\mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**}))) &= Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}})) + \lambda^2 Cov(\mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) \\ &\quad + 2\lambda Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) \\ &\geq Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}})) \end{aligned}$$

$$\Rightarrow \lambda^2 Cov(\mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) + 2\lambda Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) \geq \mathbf{0}_{q \times q}$$

Hieraus ergibt sich

$$Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) = \mathbf{0}_{q \times q}, \quad \text{da}$$

wenn $Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) < \mathbf{0}_{q \times q}$, $\exists \lambda \in \mathbb{R}$, so dass die Matrix

$$\lambda^2 Cov(\mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) + \lambda Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) < \mathbf{0}_{q \times q}.$$

ist. Wenn $Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) > \mathbf{0}_{q \times q}$, $\exists \lambda \in \mathbb{R}$, mit

$$\lambda^2 Cov(\mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) + \lambda Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) < \mathbf{0}_{q \times q}.$$

ist, d.h. $\mathbf{Y}(\tilde{\boldsymbol{\beta}})$ nicht BUP, wenn $Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*) - \mathbf{Y}(\boldsymbol{\beta}^{**})) \neq \mathbf{0}_{q \times q}$

Hinreichend: $\mathbf{Y}(\boldsymbol{\beta}^*)$ ist eine lineare erwartungstreue Vorhersage für \mathbf{Y} , dann ist $E(\mathbf{Y}(\tilde{\boldsymbol{\beta}}) - \mathbf{Y}(\boldsymbol{\beta}^*)) = \mathbf{0}_{q \times 1}$,

$$Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\tilde{\boldsymbol{\beta}}) - \mathbf{Y}(\boldsymbol{\beta}^*)) = \mathbf{0}_{q \times q}$$

und

$$\begin{aligned} Cov(\mathbf{Y} - \mathbf{Y}(\boldsymbol{\beta}^*)) &= Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}) + \mathbf{Y}(\tilde{\boldsymbol{\beta}}) - \mathbf{Y}(\boldsymbol{\beta}^*)) \\ &= Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}})) + Cov(\mathbf{Y}(\tilde{\boldsymbol{\beta}}) - \mathbf{Y}(\boldsymbol{\beta}^*)) \\ &\quad + 2Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}}), \mathbf{Y}(\tilde{\boldsymbol{\beta}}) - \mathbf{Y}(\boldsymbol{\beta}^*)) \\ &= Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}})) + Cov(\mathbf{Y}(\tilde{\boldsymbol{\beta}}) - \mathbf{Y}(\boldsymbol{\beta}^*)). \end{aligned}$$

Damit ist $Cov(\mathbf{Y} - \mathbf{Y}(\tilde{\boldsymbol{\beta}})) \leq Cov(\mathbf{Y} - \mathbf{Y}(\boldsymbol{\beta}^*))$ bzw. $\mathbf{Y}(\tilde{\boldsymbol{\beta}})$ BLUP.

Spezialfall:

Wir betrachten das Modell mit der Annahme, dass $\mathbf{V}_{xy} = \mathbf{0}_{n \times q}$ ist, dann

$V = \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{0} \\ \mathbf{0} & \mathbf{V}_{yy} \end{pmatrix}$. Aus der Theorie der linearen Modelle erhalten wir:

$$\begin{aligned} \hat{\beta}_{GM} &= \arg \min_{\beta} (\mathbf{X} - \mathbf{F}\beta)^T \mathbf{V}_{xx}^{-1} (\mathbf{X} - \mathbf{F}\beta) \\ &= (\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{X}. \\ &\quad \text{BLUE für } \beta \\ \hat{\beta}_{OLS} &= \arg \min_{\beta} (\mathbf{X} - \mathbf{F}\beta)^T (\mathbf{X} - \mathbf{F}\beta) \\ &= (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{X}. \end{aligned}$$

Statistische Eigenschaft von $\hat{\beta}_{GM}$ und $\hat{\beta}_{OLS}$:

- Erwartungswert:

$$\begin{aligned} E(\hat{\beta}_{GM}) &= E((\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{X}) \\ &= (\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F} \beta \\ &= \beta. \\ E(\hat{\beta}_{OLS}) &= E((\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{X}) \\ &= (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{F} \beta \\ &= \beta. \end{aligned}$$

- Kovarianzmatrix:

$$\begin{aligned} Cov(\hat{\beta}_{GM}) &= Cov((\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{X}) \\ &= (\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{V}_{xx} \mathbf{V}_{xx}^{-1} \mathbf{F} (\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \\ &= (\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \\ Cov(\hat{\beta}_{OLS}) &= Cov((\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{X}) \\ &= (\mathbf{F}^T \mathbf{F})^{-1} (\mathbf{F}^T \mathbf{V}_{xx} \mathbf{F}) (\mathbf{F}^T \mathbf{F})^{-1}. \end{aligned}$$

Hieraus ergibt sich

$$\begin{aligned} \mathbf{Y}(\hat{\beta}_{GM}) &= \mathbf{F}_0 (\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{X}. \\ \mathbf{Y}(\hat{\beta}_{OLS}) &= \mathbf{F}_0 (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{X} \end{aligned}$$

Statistische Eigenschaft der besten linearen Vorhersage:

- Erwartungswert und Kovarianzmatrix:

$$\begin{aligned} E(\mathbf{Y}(\hat{\beta}_{GM})) &= \mathbf{F}_0 (\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx} \mathbf{F} \beta \\ &= \mathbf{F}_0 \beta = E(\mathbf{Y}). \\ E(\mathbf{Y}(\hat{\beta}_{OLS})) &= \mathbf{F}_0 (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{F} \beta = \mathbf{F}_0 \beta = E(\mathbf{Y}). \end{aligned}$$

4. Vorhersage in Linearen Modellen

$$\begin{aligned} \text{Cov}(\mathbf{Y}(\widehat{\boldsymbol{\beta}}_{GM})) &= \mathbf{F}_0(\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}_0^T. \\ \text{Cov}(\mathbf{Y}(\widehat{\boldsymbol{\beta}}_{OLS})) &= \mathbf{F}_0(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T \mathbf{V}_{xx} \mathbf{F}(\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}_0^T. \end{aligned}$$

und

$$\text{Cov}(\mathbf{Y}(\widehat{\boldsymbol{\beta}}_{GM}), \mathbf{Y}(\widehat{\boldsymbol{\beta}}_{OLS})) = \mathbf{F}_0(\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}_0^T.$$

- Hieraus ergibt sich:

$$\begin{aligned} \text{Cov}(\mathbf{Y} - \mathbf{Y}(\widehat{\boldsymbol{\beta}}_{GM}), \mathbf{Y}(\widehat{\boldsymbol{\beta}}_{GM}) - \mathbf{Y}(\widehat{\boldsymbol{\beta}}_{OLS})) &= \text{Cov}(\mathbf{Y}(\widehat{\boldsymbol{\beta}}_{GM}), \mathbf{Y}(\widehat{\boldsymbol{\beta}}_{OLS}) - \mathbf{Y}(\widehat{\boldsymbol{\beta}}_{GM})) \\ &= \mathbf{F}_0(\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}_0^T - \mathbf{F}_0(\mathbf{F}^T \mathbf{V}_{xx}^{-1} \mathbf{F})^{-1} \mathbf{F}_0^T \\ &= \mathbf{0}. \end{aligned}$$

Durch obigen Spezialfall ergibt sich die Bemerkung:

Bemerkung 4.1 *Im Fall, dass \mathbf{X} und \mathbf{Y} unabhängig sind, hat das Lemma 4.3 die folgende Form:*

Sei $\widetilde{\boldsymbol{\beta}}$ linearer erwartungstreuer Schätzer, dann gilt:

$$\widetilde{\boldsymbol{\beta}} \text{ ist BLUP} \iff \text{Cov}(\mathbf{Y} - \mathbf{Y}(\widetilde{\boldsymbol{\beta}}), \mathbf{Y}(\boldsymbol{\beta}^*)) = \mathbf{0}, \quad \forall \boldsymbol{\beta}^* \in U_0.$$

Beispiel 4.2 (Autoregressive Kovarianz)

Wir betrachten das Modell:

$$\begin{pmatrix} X_1 \\ X_2 \\ Y \end{pmatrix} \sim N_3 \left(\begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} \boldsymbol{\beta}, \begin{pmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{pmatrix} \right), \quad \text{mit } |\rho| \leq 1 \text{ bekannt.}$$

Parameter des Modells ist: $\boldsymbol{\beta} \in \mathbb{R}$.

Modellannahme: \mathbf{X} beobachtbar, Y nicht beobachtbar. Aus der Theorie der linearen Modelle und mit Hilfe der obigen Definitionen und Theoreme sieht man, dass

- *Der Gauss Markov-Schätzer, ML-Schätzer ist:*

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_{GM} &= (\mathbf{1}_2^T \mathbf{V}_{xx}^{-1} \mathbf{1}_2)^{-1} \mathbf{1}_2^T \mathbf{V}_{xx}^{-1} \mathbf{X} \\ &= (\mathbf{1}_2^T \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix} \mathbf{1}_2)^{-1} \mathbf{1}_2^T \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix} \mathbf{X} \\ &= \frac{1}{2(1-\rho)} (1-\rho)(X_1 + X_2) \\ &= \bar{\mathbf{X}} \\ &= \widehat{\boldsymbol{\beta}}_{OLS}. \end{aligned}$$

- *Die beste lineare erwartungstreue Vorhersage ist:*

$$\begin{aligned} Y(\widehat{\boldsymbol{\beta}}_{GM}) &= \widehat{\boldsymbol{\beta}}_{GM} + \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} (\mathbf{X} - \widehat{\boldsymbol{\beta}}_{GM}) \\ &= \widehat{\boldsymbol{\beta}}_{GM} + \frac{\rho}{1-\rho^2} \begin{pmatrix} \rho & 1 \end{pmatrix} \begin{pmatrix} 1 & -\rho \\ -\rho & 1 \end{pmatrix} (\mathbf{X} - \widehat{\boldsymbol{\beta}}_{GM}) \\ &= \bar{\mathbf{X}} + \rho \begin{pmatrix} 0 & 1 \end{pmatrix} (\mathbf{X} - \bar{\mathbf{X}}) \\ &= \bar{\mathbf{X}} + \rho(\mathbf{X}_2 - \bar{\mathbf{X}}) \\ &= \frac{(1-\rho)X_1 + (1+\rho)X_2}{2}. \end{aligned}$$

- Wegen Lemma(4.3) ist

$$\begin{aligned}
\text{Cov}(Y - Y(\hat{\beta}_{GM}), X_2 - X_1) &= \text{Cov}(Y, X_2 - X_1) - \text{Cov}(Y(\hat{\beta}_{GM}), X_2 - X_1) \\
&= \rho - \rho^2 - \frac{1}{2} \text{Cov}((1 - \rho)X_1 + (1 + \rho)X_2, X_2 - X_1) \\
&= \rho - \rho^2 - \frac{1}{2}(2\rho - 2\rho^2) = 0
\end{aligned}$$

4.2. Vorhersageintervalle

Die obige Vorgehensweise lässt sich auf die Vorhersage eines Werts anwenden. Die Idee besteht darin die gemeinsame Verteilung des Zufallsvektors X und der unbeobachtbaren Zufallsvariablen Y zu betrachten. Es ist allerdings unerlässlich, neben einem Vorhersagewert stets eine Angabe über seine Qualität oder seine Präzision zu machen. So kann man beispielsweise mit einigen wenigen Beobachtungen einen Vorhersagewert ausrechnen und diesen angeben, dieser hat aufgrund seiner großen Varianz eine geringe Aussagekraft. Erst durch eine ausreichend hohe Stichprobenzahl kann eine hinreichende Präzision garantiert werden. Natürlich hängt die Präzision immer mit dem gewählten Modell und der Aufgabenstellung zusammen, so dass allein die Größe der Stichproben allein kein zuverlässiges Qualitätsmerkmal darstellt. Ein zuverlässiges und allgemeines Merkmal für die Qualität einer Vorhersage ist ein Vorhersageintervall. Für Beweise und weitere Ergebnisse verweisen wir auf die Literatur, z.B. Shao (2003, Kap.7.1.4), Rao, Shalabh, Toutenburg und Heumann (2008, Kap.6.7) und Christensen (2013).

Grundkonzept:

Definition 4.3 (*Vorhersagebereich, $(1 - \alpha)$ -Vorhersagebereich*)

(i) Eine Abbildung $\mathcal{I} : \mathbb{R}^n \rightarrow \mathcal{B}(\mathbb{R}^q)$ heißt Vorhersagebereich für \mathbf{Y} , wenn

$$\{\mathbf{x} \in \mathbb{R}^n : \mathbf{z} \in \mathcal{I}(\mathbf{x})\} \in \mathcal{B}^n, \forall \mathbf{z} \in \mathbb{R}^q.$$

(ii) Die Wahrscheinlichkeit

$$P(\mathbf{Y} \in \mathcal{I}(\mathbf{X})), \quad (4.8)$$

heißt die Überdeckungswahrscheinlichkeit des Vorhersagebereichs \mathcal{I} .

(iii) Für ein kleines Toleranzniveau $\alpha \in (0, 1)$ heißt \mathcal{I} ein $(1 - \alpha)$ -Vorhersagebereich, wenn gilt:

$$P(\mathbf{Y} \in \mathcal{I}(\mathbf{X})) \geq 1 - \alpha. \quad (4.9)$$

Beachte: Ziel ist die Konstruktion von $(1 - \alpha)$ -Vorhersagebereichen für ein vorgegebenes kleines α -Niveau (z.B. $\alpha = 0,05$), die möglichst kleine Bereiche liefert, welche die Überdeckungseigenschaft (4.9) erfüllen, das heißt die Vorhersagebereiche sollten möglichst klein sein und möglichst wenige falsche Werte überdecken.

Beispiele und Bemerkungen:

Bei einer Prognose macht man zwangsläufig einen Prognosefehler. Der Prognosefehler ist die Abweichung der Prognose von dem wahren Wert \mathbf{y} . Mit Hilfe der Eigenschaften der Parameterschätzer kann ein Intervall für den Prognosefehler hergeleitet werden. Dies ist zur Beurteilung der Prognose oft sehr hilfreich. Wir untersuchen diesen Begriff des Vorhersageintervalls in einigen Beispielen.

(i) Wir gehen wieder von folgendem Modell aus:

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim N_{n+q} \left(\begin{bmatrix} \mathbf{F}\boldsymbol{\beta} \\ \mathbf{F}_0\boldsymbol{\beta} \end{bmatrix}, \sigma^2 \begin{bmatrix} \mathbf{V}_{\mathbf{xx}} & \mathbf{V}_{\mathbf{xy}} \\ \mathbf{V}_{\mathbf{xy}}^T & \mathbf{V}_{\mathbf{yy}} \end{bmatrix} \right) \text{ mit } (\boldsymbol{\beta}, \sigma^2) \in \mathbb{R}^p \times (0, \infty).$$

Die beste lineare erwartungstreue Vorhersage für \mathbf{Y} ist gegeben durch:

$$\mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}) = \mathbf{F}_0\hat{\boldsymbol{\beta}}_{GM} + \mathbf{V}_{\mathbf{xy}}^T \mathbf{V}_{\mathbf{xx}}^{-1}(\mathbf{X} - \mathbf{F}\hat{\boldsymbol{\beta}}_{GM}).$$

Wegen $\hat{\boldsymbol{\beta}}_{GM} \sim N_p(\boldsymbol{\beta}, \sigma^2(\mathbf{F}^T \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{F})^{-1})$ gilt für die Linearkombination

$$\mathbf{Y} - \mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}) \sim N_q(0, \sigma^2 \mathbf{D}),$$

wobei

$$\mathbf{D} = \text{Cov}(\mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}) - \mathbf{Y})$$

Aus den Verteilungsergebnissen der Normalverteilung erhalten wir, dass

$$\frac{(\mathbf{Y} - \mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}))^T \mathbf{D}^{-1} (\mathbf{Y} - \mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}))}{\sigma^2} \sim \chi_q^2$$

ist. Wenn wir σ^2 durch die Schätzung $\hat{\sigma}^2$ ersetzen, ist der resultierende Ausdruck F-verteilt mit $(q, n-p)$ Freiheitsgraden und es gilt

$$P \left((\mathbf{Y} - \mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}))^T \mathbf{D}^{-1} (\mathbf{Y} - \mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM})) \leq q F_{q, n-p, 1-\alpha} S^* \right) = 1 - \alpha,$$

mit

$$S^*(\mathbf{x}) = \hat{\sigma}^2 = \frac{1}{n-p} (\mathbf{x} - \mathbf{F}\hat{\boldsymbol{\beta}}_{GM})^T \mathbf{V}_{\mathbf{xx}}^{-1} (\mathbf{x} - \mathbf{F}\hat{\boldsymbol{\beta}}_{GM}).$$

Hieraus ergibt sich der folgende Vorhersagebereich.

Satz 4.1 *Im Fall $\text{Rang}(\mathbf{F}_0) = q$ ist ein $(1-\alpha)$ -Vorhersageellipsoid für \mathbf{Y} gegeben durch:*

$$\mathcal{I}(\mathbf{x}) = \left\{ \mathbf{y} \in \mathbb{R}^q : \frac{1}{q} \frac{(\mathbf{y} - \mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}))^T \mathbf{D}^{-1} (\mathbf{y} - \mathbf{Y}(\hat{\boldsymbol{\beta}}_{GM}))}{S^*(\mathbf{x})} \leq F_{q, n-p, 1-\alpha} \right\}. \quad (4.10)$$

- (ii) Seien $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ stochastisch unabhängig und identisch $N_p(\boldsymbol{\beta}, \mathbf{V})$ -verteilte \mathbb{R}^p -wertige Zufallsvariablen, $(\boldsymbol{\beta}, \mathbf{V}) \in \mathbb{R}^p \times PD(p)$ ist der Parameter. Kurz

$$\left((\mathbb{R}^n)^p, \mathcal{B}^{n \times p}, (\otimes_1^n N(\boldsymbol{\beta}, \mathbf{V}))_{(\boldsymbol{\beta}, \mathbf{V}) \in \mathbb{R}^p \times PD(p)} \right). \quad (4.11)$$

Sei \mathbf{X}_{n+1} ein weiterer Zufallsvektor aus dieser Verteilung, nicht beobachtbar, unabhängig von den ersten n . Wir fassen die Zufallsvariablen zu einer $\mathbb{R}^{n \times p}$ -wertigen Zufallsvariablen zusammen:

$$\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^T.$$

Das Modell schreibt sich dann so:

$$\begin{pmatrix} \text{Vec}(\mathbf{X}) \\ \mathbf{X}_{n+1} \end{pmatrix} \sim N_{np+p} \left(\begin{bmatrix} \mathbf{1}_n \otimes \boldsymbol{\beta} \\ \boldsymbol{\beta} \end{bmatrix}, \sigma^2 \begin{bmatrix} I_n \otimes \mathbf{V} & 0 \\ 0 & \mathbf{V} \end{bmatrix} \right).$$

Die optimale Prognose für eine neue (zukünftige) Beobachtung \mathbf{X}_{n+1} ist gegeben durch:

$$\begin{aligned} \widehat{\mathbf{X}}_{n+1} &= \widehat{\boldsymbol{\beta}} \\ &= \mathbf{X}^T \frac{1}{n} \mathbf{1}_n \\ &= \bar{\mathbf{X}}. \end{aligned}$$

Neben der Punktschätzung der Prognose ist man in der Regel auch an einem Vorhersageintervall interessiert. Ein Vorhersageintervall für \mathbf{X}_{n+1} lässt sich folgendermaßen konstruieren. Wegen

$$\begin{aligned} \mathbf{D} = \text{Cov}(\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1}) &= \text{Cov}(\mathbf{X}_{n+1}) + \text{Cov}(\widehat{\mathbf{X}}_{n+1}) - \text{Cov}(\mathbf{X}_{n+1}, \widehat{\mathbf{X}}_{n+1}) \\ &\quad - \text{Cov}(\widehat{\mathbf{X}}_{n+1}, \mathbf{X}_{n+1}) \\ &= \mathbf{V} + \frac{\mathbf{V}}{n}. \end{aligned}$$

gilt

$$\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1} \sim N_p(0, \mathbf{D}).$$

Kleines Hilfsresultat: Seber(1984, S. 30)

Wenn $\mathbf{Y} \sim N_p(0, \mathbf{V})$, $\mathbf{A} \sim W_p(m, \mathbf{V})$ und \mathbf{Y}, \mathbf{A} unabhängig sind, dann gilt

$$\frac{m-p+1}{p} \mathbf{Y}^T \mathbf{A}^{-1} \mathbf{Y} \sim F_{p, m-p+1} \quad (4.12)$$

Wir betrachten die Statistik:

$$\begin{aligned} \mathbf{S}(\mathbf{x}) &= \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \\ &= \mathbf{x}^T \mathbf{x} - n \bar{\mathbf{x}} \bar{\mathbf{x}}^T. \end{aligned}$$

4. Vorhersage in Linearen Modellen

Durch geeignetes Standardisieren erhalten wir:

$$\mathbf{Q} = \sqrt{\frac{n}{n+1}} \mathbf{D}^{-1/2} (\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1}) \sim N_p(0, I_p)$$

$$\mathbf{B} = \mathbf{V}^{-1/2} \mathbf{S} \mathbf{V}^{-1/2} \sim W_p(n-1, I_p).$$

Im Hinblick, dass \mathbf{B} und \mathbf{Q} unabhängig sind, und unter Verwendung der Darstellung (4.12) erhalten wir

$$\frac{n-p}{p} \mathbf{Q}^T \mathbf{B}^{-1} \mathbf{Q} \sim F_{p, n-p}.$$

Wenn wir \mathbf{Q} durch $\sqrt{\frac{n}{n+1}} \mathbf{D}^{-1/2} (\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1})$ ersetzen ist der resultierende Ausdruck

$$\frac{(n-p)n}{(n+1)p} (\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1})^T \mathbf{S}^{-1} (\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1}) \sim F_{p, n-p}$$

und es gilt

$$P \left(\frac{(n-p)n}{(n+1)p} (\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1})^T \mathbf{S}^{-1} (\mathbf{X}_{n+1} - \widehat{\mathbf{X}}_{n+1}) \leq F_{p, n-p, 1-\alpha} \right) = 1 - \alpha.$$

Somit erhält man

$$\mathcal{I}(\mathbf{x}) = \left\{ \mathbf{x}_{n+1} \in \mathbb{R}^p : \frac{(n-p)n}{(n+1)p} (\mathbf{x}_{n+1} - \bar{\mathbf{x}})^T \mathbf{S}^{-1}(\mathbf{x}) (\mathbf{x}_{n+1} - \bar{\mathbf{x}}) \leq F_{p, n-p, 1-\alpha} (1 - \alpha) \right\}, \quad (4.13)$$

als Vorhersageintervall für \mathbf{X}_{n+1} zum Niveau $1 - \alpha$.

4.3. Weitere simultane Vorhersageintervalle unter Normalverteilungsannahme

Dieser Abschnitt stellt zunächst Vorhersageintervalle im ein- und mehrdimensionalen Fall vor und behandelt danach das simultane Vorhersageintervall nach den Ansätzen von Bonferroni und Scheffé. Für Beweise und weitere Eigenschaften verweisen wir auf die Literatur, z.B. Hocking (2005, Kap. 18.2) und Christensen(2011).

Wir betrachten das Modell 1.3.b mit \mathbb{R}^{p+q} -wertigen, normal-verteilten Zufallsvariablen,

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim N_{p+q} \left(\begin{pmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{pmatrix}, \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \right),$$

wobei

- $\mathbf{X} \sim N_p(\boldsymbol{\mu}_x, \mathbf{V}_{xx})$ (beobachtbar)
- \mathbf{Y} eine \mathbb{R}^q -wertige Zufallsvariable (nicht beobachtbar)

Als BLUP-Vorhersage $\widehat{\mathbf{Y}}$ für \mathbf{Y} ergibt sich

$$\begin{aligned}\widehat{\mathbf{Y}} &= \arg \min_{g \in SP} E[\mathbf{Y} - g(\mathbf{X})]^2 \\ &= E(\mathbf{Y}|\mathbf{X}) \\ &= \boldsymbol{\mu}_y + \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1}(\mathbf{X} - \boldsymbol{\mu}_x) \quad (\text{nach 1.12})\end{aligned}$$

Verteilung der Residuen $\mathbf{Y} - \widehat{\mathbf{Y}}$

Aus den Verteilungsergebnissen für (mehrdimensionale) normalverteilte Zufallsvariablen erhalten wir, dass

1. $\mathbf{Y} - \widehat{\mathbf{Y}} \sim N_q(0, \mathbf{C})$ das heißt $\mathbf{C}^{-1/2}(\widehat{\mathbf{Y}} - \mathbf{Y}) \sim N_q(0, I_q)$, wobei

$$\mathbf{C} = \text{Cov}(\mathbf{Y} - \widehat{\mathbf{Y}}) = \mathbf{V}_{yy} - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \quad (4.14)$$

2. Hieraus ergibt sich, dass

$$(\mathbf{Y} - \widehat{\mathbf{Y}})^T \mathbf{C}^{-1} (\mathbf{Y} - \widehat{\mathbf{Y}}) \sim \chi_q^2$$

3. $Y_i - \widehat{Y}_i \sim N(0, C_{ii}) \Rightarrow \frac{Y_i - \widehat{Y}_i}{\sqrt{C_{ii}}} \sim N(0, 1)$ mit $C_{ii} = \mathbf{e}_i^T \mathbf{C} \mathbf{e}_i$ und \mathbf{e}_i ist i -ter Einheitsvektor.

Prognoseintervall für einzelne Koeffizienten Y_i

- Ein $(1 - \alpha)$ Vorhersageintervall für Y_i ist gegeben durch:

$$\mathcal{I}^{(i)}(\mathbf{x}) = [\widehat{y}_i(\mathbf{x}) - \sqrt{C_{ii} z_{1-\frac{\alpha}{2}}}, \widehat{y}_i(\mathbf{x}) + \sqrt{C_{ii} z_{1-\frac{\alpha}{2}}}],$$

wobei $z_{1-\frac{\alpha}{2}} = \Phi^{-1}(1 - \frac{\alpha}{2})$ das $(1 - \frac{\alpha}{2})$ -Quantil von $N(0, 1)$ ist.

- Die Länge des zufälligen Intervalls ist bekannt (hängt nicht von \widehat{Y}_i ab) gleich

$$2\sqrt{C_{ii} z_{1-\frac{\alpha}{2}}}$$

Simultanes Prognoseintervall (Prognosebereich) für \mathbf{Y} :

Hocking(2005, Kap.18.1.1.2)

Oft interessieren nicht nur Vorhersageintervalle für eine einzelne Komponente Y_i , sondern auch Intervalle für mehrere Kombinationen $\mathbf{e}_i^T \mathbf{Y}$, $i = 1, \dots, r$, z.B. alle Komponenten $Y_i, i = 1, \dots, q$. Ein $(1 - \alpha)$ Vorhersagebereich \mathcal{I} für \mathbf{Y} heißt ein $(1 - \alpha)$ simultanes Vorhersageintervall für $Y_i, i = 1, \dots, q$, wenn \mathcal{I} von der folgenden Form ist:

$$\mathcal{I}(\mathbf{x}) = \times_{i=1}^q \mathcal{I}^{(i)}(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^p,$$

4. Vorhersage in Linearen Modellen

mit Intervall $\mathcal{I}^{(i)}(\mathbf{x}) \subset \mathbb{R}$ für alle $i = 1, \dots, q$ und $\mathbf{x} \in \mathbb{R}^p$, d.h. $\mathcal{I}(\mathbf{x})$ ein achsenparalleler Quader in \mathbb{R}^q für jedes \mathbf{x} .

Ein simultanes $(1 - \alpha)$ Vorhersageintervall für \mathbf{Y} nach der Bonferroni-Methode ist gegeben durch:

$$\mathcal{I}_{Bonf}(\mathbf{x}) = \times_{i=1}^q \left[\mathbf{e}_i^T \hat{\mathbf{y}}(\mathbf{x}) - \sqrt{C_{ii} z_{1-\frac{\alpha}{2q}}}, \mathbf{e}_i^T \hat{\mathbf{y}}(\mathbf{x}) + \sqrt{C_{ii} z_{1-\frac{\alpha}{2q}}} \right], \quad (4.15)$$

wobei $z_{1-\frac{\alpha}{2q}}$ das $(1 - \frac{\alpha}{2q})$ -Quantil von $N(0, 1)$ ist.

Ein anderes $(1 - \alpha)$ simultanes Vorhersageintervall für die \mathbb{R}^q -wertige Zufallsvariable \mathbf{Y} ist das Vorhersageellipsoid nach Scheffé.

Theorem 4.1 (Vorhersageellipsoid nach Scheffé)

Ein $(1 - \alpha)$ Vorhersageellipsoid für \mathbf{Y} ist gegeben durch:

$$\mathcal{I}_{Sch}(\mathbf{x}) = \{ \mathbf{y} \in \mathbb{R}^q : (\mathbf{y} - \hat{\mathbf{y}}(\mathbf{x}))^T \mathbf{C}^{-1} (\mathbf{y} - \hat{\mathbf{y}}(\mathbf{x})) \leq \chi_{q,1-\alpha}^2 \}, \quad (4.16)$$

wobei $\mathcal{I}(\mathbf{x})$ ein Ellipsoid in \mathbb{R}^q mit dem Mittelpunkt $\hat{\mathbf{y}}(\mathbf{x})$ ist.

Ein ellipsenförmiger Bereich als Vorhersageintervall für die \mathbb{R}^q -wertige Vorhersage ist in den Anwendungen nicht so gut interpretierbar wie ein Achsenparalleler Quader.

Herleitung von marginalen Vorhersageintervallen nach Scheffé:

Lemma 4.4 (andere Interpretation des Vorhersageellipsoids)

Für gegebenes $\alpha \in (0, 1)$ gilt:

$$P\left(|Y_i - \hat{Y}_i| \leq \sqrt{C_{ii} \chi_{q,1-\alpha}^2}\right) = 1 - \alpha$$

Beweisidee:

Theorem 4.1 sagt, dass

$$P\left((\mathbf{Y} - \hat{\mathbf{Y}})^T \mathbf{C}^{-1} (\mathbf{Y} - \hat{\mathbf{Y}}) \leq \chi_{q,1-\alpha}^2\right) = 1 - \alpha. \quad (4.17)$$

Mit A.4 können wir schreiben:

$$(\mathbf{y} - \hat{\mathbf{y}})^T \mathbf{C}^{-1} (\mathbf{y} - \hat{\mathbf{y}}) = \max_{\mathbf{a} \in \mathbb{R}^q} \frac{\left(\mathbf{a}^T (\mathbf{y} - \hat{\mathbf{y}})\right)^2}{\mathbf{a}^T \mathbf{C} \mathbf{a}} \quad (4.18)$$

Folgerung 4.2 (marginale Vorhersageintervall) Hocking (2005, Kap.18.1.1.3)

Das marginale $(1 - \alpha)$ -Vorhersageintervall für \mathbf{Y} nach Scheffé (Scheffé Rechteck Abk.: Schr) ist gegeben durch:

$$\mathcal{I}_{Schr}(\mathbf{x}) = \times_{i=1}^q \left[\mathbf{e}_i^T \hat{\mathbf{y}}(\mathbf{x}) - \sqrt{C_{ii} \chi_{q,1-\alpha}^2}, \mathbf{e}_i^T \hat{\mathbf{y}}(\mathbf{x}) + \sqrt{C_{ii} \chi_{q,1-\alpha}^2} \right]. \quad (4.19)$$

wobei $\chi_{q,1-\alpha}^2$ das $(1 - \alpha)$ -Quantil der χ_q^2 -Verteilung ist.

4.3. Weitere simultane Vorhersageintervalle unter Normalverteilungsannahme

Beachte: Die Intervall-Längen der Bonferroni-Methode und der Scheffé-Methode sind unterschiedlich. Im vorliegenden Fall läuft es auf den Vergleich der Quantile

$$z_{1-\frac{\alpha}{2q}} \quad \text{und} \quad \sqrt{\chi_{q,1-\alpha}^2}$$

hinaus. Zum Beispiel:

$\alpha = .01$			$\alpha = .05$		
q	$z_{1-\frac{\alpha}{2q}}$	$\sqrt{\chi_{q,1-\alpha}^2}$	q	$z_{1-\frac{\alpha}{2q}}$	$\sqrt{\chi_{q,1-\alpha}^2}$
2	2.807034	3.034854	2	2.241403	2.447747
3	2.935199	3.368214	3	2.39398	2.795483
4	3.023341	3.643721	4	2.497705	3.080216
5	3.090232	3.884105	5	2.575829	3.327236
10	3.290527	4.817598	10	2.807034	4.278672
30	3.587915	7.133876	30	3.14398	6.616115

In den obigen Beispielen gilt:

$$I_{Bonf}(\mathbf{x}) \subseteq I_{Schr}(\mathbf{x}), \quad \forall \mathbf{x} \in \mathbb{R}^p$$

Daher liegt die Vermutung nahe, dass die Bonferroni-Methode besser ist als die Scheffé-Methode. Es bleibt Ansatzpunkt für weitere Untersuchungen. Eine Fragestellung für die Zukunft besteht darin, ob es einen Fall gibt, so dass $I_{Schr}(\mathbf{x}) \subseteq I_{Bonf}(\mathbf{x})$ ist.

5. Partielle kleinste Quadrate

In vielen Anwendungsfeldern der Statistik wie Medizin oder Ökonomie wird eine große Anzahl von Variablen erhoben, auch weil man nicht von vornherein weiß, welches die eigentlich wichtigen sind. Das macht die Entdeckung von Strukturen in den Daten schwierig, zumal die Variablen oft stark untereinander korrelieren, mithin gleichartige Information enthalten, ist es selten ratsam, einfach Variablen wegzulassen. Als Alternative können wir geeignete Linearkombinationen der Variablen suchen, die in geringerer Zahl möglichst alle relevanten Informationen des Datensatzes enthalten.

Für die Hauptkomponentenanalyse als datenanalytische Technik und Dimensionsreduktion kann man sagen, dass die Hauptkomponentenanalyse eine orthogonale Transformation im m -dimensionalen Raum der Originalvariablen in eine neue Variablenmenge ist. Als Dimensionsreduktion hat die Hauptkomponentenanalyse das Ziel, eine möglichst geringe Anzahl von Linearkombinationen der beobachteten Variablen zu finden, welche die Daten, möglichst gut repräsentieren. In manchen Anwendungen sind die resultierenden, als Hauptkomponenten bezeichneten Variablen die Größen, die von Interesse sind.

Die Hauptkomponentenanalyse besitzt jedoch einige Nachteile. Die Ergebnisse sind abhängig von der Skalierung und daher nicht eindeutig. Die Hauptkomponenten sind schwierig zu interpretieren. Die Interpretation ist zumeist subjektiv. Auch die Anzahl der auszuwählenden wesentlichen Hauptkomponenten ist nicht eindeutig bestimmt.

Eine weitere Möglichkeit zur Dimensionsreduktion bietet die sogenannte Methode der partiellen kleinsten Quadrate (engl. partial least squares, PLS), die ursprünglich von Wold (1966) eingeführt wurde. Die PLS-Methode hat gewisse Ähnlichkeiten mit der Hauptkomponentenanalyse PCA (engl. Principal component analysis). Bei der PCA-Methode werden ebenfalls zueinander orthogonale neue Variablen als Linearkombination aus den ursprünglichen Variablen bestimmt. Im Unterschied zur PLS werden bei der PCA (als Ansatz für Dimensionsreduktion und Vorhersage) nicht die Kovarianz zwischen Y und $\mathbf{a}^T \mathbf{X}$, sondern nur die Varianz von $\mathbf{a}^T \mathbf{X}$ maximiert (siehe Puntanen und Styan und Isotalo (2013, Kap.9.4)). Sie wird in verschiedenen Anwendungsgebieten zur Regression oder Diskriminanzanalyse und Vorhersage eingesetzt, insbesondere für die Analyse hochdimensionaler Daten. Besonders verbreitet ist die Anwendung der PLS-Methode in der Chemometrie zur spektrometrischen Kalibrierung. Prädiktoren sind dabei die Emissionsintensitäten zu den verschiedenen Frequenzen eines gewissen Spektrums, abhängige Variablen sind die Anteile von bestimmten Stoffen in einem Stoffgemisch (vgl. Grüning (2005, Kap.4.2.1)).

Im Folgenden werden wir zunächst die theoretischen Grundlagen der Methode der Partiiellen kleinsten Quadrate studieren, um das mathematische Modell und die ursprünglichen Ideen der Methode aufzubauen, die zur Entwicklung der PLS-Methode führen, um einige wichtige Eigenschaften der PLS-Methode herleiten zu können. Wir stellen die

5. Partielle kleinste Quadrate

Methoden der Hauptkomponentenanalyse und die Methode der kanonischen Korrelation als Minimierung des mittleren quadratischen Vorhersagefehlers vor. Zum Schluss betrachten wir Beziehungen zwischen den Methoden der Regression, Hauptkomponentenanalyse, kanonischen Korrelation und Partiiellen kleinsten Quadrate.

Seien X und Y quadratintegrierbare reelle Zufallsvariablen, dann betrachten wir als Motivation das Optimierungsproblem:

$$\min_{a,b \in \mathbb{R}} E[(Y - b - aX)^2]. \quad (5.1)$$

Dieses besitzt die Lösung:

$$a^* = \frac{Cov(X, Y)}{Var(X)}, \quad b^* = E(Y) - a^*E(X). \quad (5.2)$$

Es gilt $Cov(a^*X, Y) = Var(a^*X) = \frac{[Cov(X, Y)]^2}{Var(X)}$. Der Minimalwert M^* in (5.1) heißt Residualvarianz und ergibt sich zu :

$$\begin{aligned} M^* &= E[(Y - a^*X - b^*)^2] \\ &= Var(Y - a^*X). \\ &= Var(Y) - Cov(a^*X, Y). \\ &= Var(Y)(1 - Corr^2(X, Y)), \end{aligned}$$

5.1. Motivation

Wir betrachten in diesem Abschnitt das Modell 1.3.b mit $q = 1$, und treffen für das Modell folgende zusätzliche Annahmen

1. Erwartungswert ist null $E(\mathbf{X}) = 0$
2. Y ist nicht beobachtbar.

Kurz:

$$\begin{aligned} \begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix} &\sim N_{p+1}(\mathbf{0}, \mathbf{V}), \text{ mit} \\ \mathbf{V} &= \begin{pmatrix} \mathbf{V}_{\mathbf{xx}} & \mathbf{V}_{\mathbf{xy}} \\ \mathbf{V}_{\mathbf{yx}} & V_{yy} \end{pmatrix} \in PD(p+1) \\ &Y \text{ nicht beobachtbar} \end{aligned}$$

Wir betrachten das folgende Problem der Prognose: basierend auf einem Beobachtungswert der \mathbb{R}^p -wertigen Zufallsvariable \mathbf{X} , möchte man eine beste lineare Vorhersage von Y möglichst im Sinne der Minimierung des mittleren quadratischen Fehlers herleiten. Also ist X eine beobachtete Zufallsvariable, die die Variable Y gut vorhersagen soll. Mit anderen Worten hat unsere Aufgabe häufig das Ziel, eine lineare Transformation $\mathbf{a}^T \mathbf{X}$, mit $\mathbf{a}^T \mathbf{a} = 1$ so bestimmt werden, dass $E[(Y - \mathbf{a}^T \mathbf{X})^2]$ minimal wird.

(I) Aus den Eigenschaften der bedingten Erwartung folgt:

$$\min_{g \in SP} E[(Y - g(\mathbf{X}))^2] = E[(Y - E(Y|\mathbf{X}))^2] \quad (5.3)$$

$$= \text{Var}(Y - \mathbf{V}_{yx} \mathbf{V}_{xx}^{-1} \mathbf{X}). \quad (5.4)$$

Aus der Definition der bedingten Erwartung der Normalverteilung (1.13) bemerken wir, dass die optimale Lösung von $\min_{g \in SP} \text{Var}(Y - g(\mathbf{X}))$ mit der optimalen Lösung von $\min_{l \in LP} \text{Var}(Y - l(\mathbf{X}))$ unter Normalverteilungsannahme übereinstimmt.

(II) Aus den Verteilungsergebnissen für (mehrdimensionale) normalverteilte Zufallsvariablen wissen wir, dass

$$\begin{pmatrix} \mathbf{a}^T \mathbf{X} \\ Y \end{pmatrix} \sim N_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \mathbf{a}^T \mathbf{V}_{xx} \mathbf{a} & \mathbf{a}^T \mathbf{V}_{xy} \\ \mathbf{V}_{yx} \mathbf{a} & V_{yy} \end{pmatrix} \right)$$

ist. (siehe Fujikoshi, Ulyanov und Shimizu (2010, Aufgabe 11.7))

- Die beste lineare erwartungstreue Vorhersage für Y unter $\mathbf{a}^T \mathbf{X}$ ist gegeben durch:

$$\hat{Y}_{\mathbf{a}} = E(Y | \mathbf{a}^T \mathbf{X}) \quad (5.5)$$

$$= \text{Cov}(Y, \mathbf{a}^T \mathbf{X}) [\text{Var}(\mathbf{a}^T \mathbf{X})]^{-1} \mathbf{a}^T \mathbf{X} \quad (\text{nach (1.13)}) \quad (5.6)$$

$$= \frac{\mathbf{V}_{yx} \mathbf{a} \mathbf{a}^T}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \mathbf{X} \quad (5.7)$$

- Der mittlere quadratische Vorhersagefehler von $\hat{Y}_{\mathbf{a}}$ ist gegeben durch:

$$\begin{aligned} R(\mathbf{a}) &= E(Y - \hat{Y}_{\mathbf{a}})^2 \\ &= \text{Var}(Y - \hat{Y}_{\mathbf{a}}) \\ &= \text{Var}(Y) - \text{Var}(\hat{Y}_{\mathbf{a}}) \\ &= V_{yy} - \frac{\mathbf{a}^T \mathbf{V}_{xy} \mathbf{V}_{xy}^T \mathbf{a}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \quad (\text{nach Folgerung 3.5}) \end{aligned}$$

Gesucht wird die Linearkombination

$$\mathbf{a}^T \mathbf{X} = a_1 X_1 + \dots + a_p X_p, \quad \mathbf{a} \neq 0$$

für die $R(\mathbf{a})$ minimum ist. Minimierung von $R(\mathbf{a})$ ist äquivalent zu Maximierung von $\frac{\mathbf{a}^T \mathbf{V}_{xy} \mathbf{V}_{xy}^T \mathbf{a}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}}$, und die optimale Lösung ist :

$$\mathbf{a}^* = e.v_{max}(\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{V}_{xy}^T) \quad (\text{nach Satz A.3}), \quad (5.8)$$

wobei \mathbf{a}^* der Eigenvektor ist, der zum λ_{\max} der Matrix $\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{V}_{xy}^T$ gehört. Nach Satz A.4 erhalten wir folgende Lösung:

$$\mathbf{a}^* = \frac{\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy}}{\mathbf{V}_{xy}^T \mathbf{V}_{xx}^{-1} \mathbf{V}_{xy}}$$

5. Partielle kleinste Quadrate

- Es gilt also:

$$\begin{aligned}
 R(\mathbf{a}) &= V_{yy} - \frac{\mathbf{V}_{yx}\mathbf{a}\mathbf{a}^T\mathbf{V}_{xy}}{\mathbf{a}^T\mathbf{V}_{xx}\mathbf{a}} \\
 &= \text{Var}(Y) - \text{Cov}(Y, \hat{Y}_{\mathbf{a}}) \\
 &= \text{Var}(Y) - \text{Var}(\hat{Y}_{\mathbf{a}}) \\
 &= \text{Var}(Y)(1 - \text{Corr}^2(\hat{Y}_{\mathbf{a}}, Y)).
 \end{aligned}$$

Beachte

- Hieraus ergibt sich, dass

$$\begin{aligned}
 \min_{\mathbf{a}} R(\mathbf{a}) &= \min_{\mathbf{a}} \text{Cov}(Y - \hat{Y}_{\mathbf{a}}, Y) \\
 &= \min_{\mathbf{a}} [\text{Var}(Y) - \text{Var}(\hat{Y}_{\mathbf{a}})] \\
 &= \min_{\mathbf{a}} \text{Var}(Y)(1 - \text{Corr}^2(\hat{Y}_{\mathbf{a}}, Y)).
 \end{aligned}$$

sind.

- Eine äquivalente Darstellung lautet:

$$\min_{\mathbf{a} \in \mathbb{R}^p} R(\mathbf{a}) \Leftrightarrow \max_{\mathbf{a}} \text{Cov}(\hat{Y}_{\mathbf{a}}, Y) \quad (5.9)$$

$$\Leftrightarrow \max_{\mathbf{a}} \text{Var}(\hat{Y}_{\mathbf{a}}) \quad (5.10)$$

$$\Leftrightarrow \max_{\mathbf{a}} \text{Corr}^2(\hat{Y}_{\mathbf{a}}, Y). \quad (5.11)$$

- Hieraus ergibt sich

$$\min_{\mathbf{a} \in \mathbb{R}^p} \text{Var}(Y - \mathbf{a}^T \mathbf{X}) = \min_{\mathbf{a} \in \mathbb{R}^p} [\min_{h \in SP_1} \text{Var}(Y - h(\mathbf{a}^T \mathbf{X}))] \quad (5.12)$$

$$= \min_{\mathbf{a} \in \mathbb{R}^p} \text{Var}(Y - E(Y | \mathbf{a}^T \mathbf{X})), \quad (5.13)$$

wobei $SP_1 = \{g | g : \mathbb{R} \rightarrow \mathbb{R}, g \text{ messbar}, E(g(Z)^2) < \infty\}$ ist. In diesem Fall ist dazu äquivalent

$$\min_{\mathbf{a} \in \mathbb{R}^p} \text{Var}(Y - \mathbf{a}^T \mathbf{X}) = \min_{\mathbf{a} \in \mathbb{R}^p} [\min_{h \in L_1} \text{Var}(Y - h(\mathbf{a}^T \mathbf{X}))] \quad (5.14)$$

$$= \min_{\mathbf{a} \in \mathbb{R}^p} \text{Var}(Y - BLUP(Y, \mathbf{a}^T \mathbf{X})), \quad (5.15)$$

$$L_1 = \{g | g : \mathbb{R} \rightarrow \mathbb{R}, g \text{ messbar}, g(z) = cz, c \in \mathbb{R}\}.$$

Dieses Optimierungsproblem für $R(\mathbf{a})$ werden wir in diesem Kapitel und in den anschließenden Aufgaben vertiefen.

5.2. Modellspezifikation und Problemstellung

Wir betrachten das Modell 1.3.a mit $q = 1$, unter der zusätzlichen Annahme, dass Y nicht beobachtbar ist. Kurz:

$$\begin{aligned} & \begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix} \text{ eine } \mathbb{R}^{p+1} \text{ wertige ZV, mit} \\ & \boldsymbol{\mu} \in \mathbb{R}^{p+1}, \mathbf{V} \in PD(p+1) \\ & Y \text{ nicht beobachtbar.} \end{aligned}$$

Wir hatten im Kapitel 3 die Aufgabe, basierend auf den Beobachtungswerten x_1, \dots, x_p der \mathbb{R}^p -wertigen Zufallsvariable \mathbf{X} den Wert der \mathbb{R} -wertigen Zufallsvariable Y (im Sinne des mittleren quadratischen Fehlers) zu prognostizieren, untersucht.

$$\min_{g \in SP} E[(Y - g(\mathbf{X}))^2]. \quad (5.16)$$

Die optimale Lösung von (5.16) war $m(\mathbf{X}) = E(Y|\mathbf{X})$. Zur Bestimmung von $E(Y|\mathbf{X})$ benötigen wir die gemeinsame Verteilung von Y und \mathbf{X} . Außerdem ist $E(Y|\mathbf{X})$ nicht einfach zu bestimmen. Aus diesem Grund beschränken wir SP auf die Klasse LP .

Nachteile:

Da $LP \subseteq SP$ gilt, dass

$$\begin{aligned} \min_{g \in SP} E[(Y - g(\mathbf{X}))^2] &\leq \min_{l \in LP} E[(Y - l(\mathbf{X}))^2] \\ \max_{g \in SP} \text{Corr}^2(Y, g(\mathbf{X})) &\geq \max_{l \in LP} \text{Corr}^2(Y, l(\mathbf{X})) \end{aligned}$$

Vorteile:

1. Die Berechnung von $\min_{l \in LP} E[(Y - l(\mathbf{X}))^2]$ ist einfach.
2. $BLUP(Y, \mathbf{X})$ hängt nur von den ersten und zweiten Momenten ab
3. Unter Normalverteilungsannahme gilt:

$$\min_{g \in SP} E[(Y - g(\mathbf{X}))^2] = \min_{l \in LP} E[(Y - l(\mathbf{X}))^2].$$

Die Ausdrücke (5.9), (5.10) und (5.11) führen zu den folgenden Begriffsbildungen:

(i) Vorhersage durch Maximierung der Korrelation:

Es wurde bereits im Kapitel 3 gezeigt, dass

$$\begin{aligned} \arg \max_{g \in SP} \text{Corr}^2(Y, g(\mathbf{X})) &= E(Y|\mathbf{X}) \\ \arg \max_{\mathbf{a} \in \mathbb{R}^p} \text{Corr}^2(Y, \mathbf{a}^T \mathbf{X}) &= \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}}. \end{aligned}$$

(ii) Vorhersage durch Maximierung der Varianz oder der Kovarianz

Dazu betrachtet man

$$\text{Var}(\mathbf{a}^T \mathbf{X}) = \mathbf{a}^T \mathbf{V}_{\mathbf{xx}} \mathbf{a} \quad \text{und} \quad \text{Cov}(\mathbf{a}^T \mathbf{X}, Y) = \mathbf{a}^T \mathbf{V}_{\mathbf{xy}}. \quad (5.17)$$

Diese Vorhersagen sind bestimmt durch Maximierung der Varianz der Linearkombination $\mathbf{a}^T \mathbf{X}$ oder durch Maximierung der Kovarianz zwischen $\mathbf{a}^T \mathbf{X}$ und Y . Die Maximierungsprobleme sind nicht eindeutig lösbar auf $\mathbf{a} \in \mathbb{R}^p$. Der Vektor \mathbf{a} ist damit nur bis auf eine Konstante eindeutig bestimmt und eine zusätzliche Normierungsbedingung ist notwendig um Eindeutigkeit zu erreichen, wie zum Beispiel in der kanonischen Korrelation $\text{Var}(\mathbf{a}^T \mathbf{X}) = \mathbf{a}^T \mathbf{V}_{\mathbf{xx}} \mathbf{a} = 1$. Die Lösung von

$$\begin{aligned} & \max_{\mathbf{a}} \text{Cov}(\mathbf{a}^T \mathbf{X}, Y) \\ & \text{unter der Nebenbedingung } \mathbf{a}^T \mathbf{V}_{\mathbf{xx}} \mathbf{a} = 1, \end{aligned}$$

ist $\sqrt{\mathbf{V}_{\mathbf{xy}}^T \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}}}$. Anstatt die Bedingung $\mathbf{a}^T \mathbf{V}_{\mathbf{xx}} \mathbf{a} = 1$ zu fordern, wählen wir in diesem Kapitel die Nebenbedingung $\mathbf{a}^T \mathbf{a} = 1$ (in der Praxis wie z.B. Data Mining). Die Maximierungsprobleme stellen sich damit in der Form

$$\begin{aligned} & \max_{\mathbf{a}} \text{Var}(\mathbf{a}^T \mathbf{X}) \quad \text{oder} \quad \max_{\mathbf{a}} \text{Cov}(\mathbf{a}^T \mathbf{X}, Y) \\ & \text{unter der Nebenbedingung } \mathbf{a}^T \mathbf{a} = 1. \end{aligned}$$

In diesem Kapitel und in den anschließenden Aufgaben werden wir das Optimierungsproblem

$$\min_{g \in LP_1} E(Y - g(\mathbf{X}))^2 \quad (5.18)$$

vertiefen, wobei die Menge

$$LP_1 = \{g | g : \mathbb{R}^p \rightarrow \mathbb{R}, g(\mathbf{x}) = \mu_y + \mathbf{a}^T (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}}), \mathbf{a} \in \mathbb{R}^p \text{ mit } \mathbf{a}^T \mathbf{a} = 1\} \subseteq LP$$

kompakt ist. Durch die Maximierung von Kovarianz-Varianz lassen sich verschiedene Vorhersagen für Y unter \mathbf{X} in der Klasse LP_1 bestimmen. Ziel ist es, die Qualität einer Vorhersage für Y zu messen, wenn der Prognosefehler $Y - \hat{Y}$ nicht orthogonal zu allen in die Prognosefunktion eingehenden Variablen X_1, \dots, X_p steht.

5.3. Die Vorhersage auf LP_1

5.3.1. Vorhersage durch Maximierung der Korrelation

Wir suchen einen Koeffizientenvektor $\mathbf{a} \in \mathbb{R}^p \setminus \{0\}$ mit $\mathbf{a}^T \mathbf{a} = 1$, wobei die Korrelation zwischen Y und $\mathbf{a}^T \mathbf{X}$ maximal ist.

$$\begin{aligned} \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Corr}(Y, \mathbf{a}^T \mathbf{X}) &= \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Corr}(Y, \mathbf{a}^T \mathbf{V}_{\mathbf{xx}}^{-1/2} \mathbf{X}), \quad \text{mit } \mathbf{V}_{\mathbf{xx}}^{-1/2} \mathbf{V}_{\mathbf{xx}}^{-1/2} = \mathbf{V}_{\mathbf{xx}}^{-1} \\ &\Leftrightarrow \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Cov}(\mathbf{a}^T \mathbf{V}_{\mathbf{xx}}^{-1/2} \mathbf{X}, Y) \\ &\Leftrightarrow \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \mathbf{a}^T \mathbf{V}_{\mathbf{xx}}^{-1/2} \mathbf{V}_{\mathbf{xy}} \\ &= \sqrt{\mathbf{V}_{\mathbf{yx}} \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{yx}}} \quad (\text{Satz A.1}), \end{aligned}$$

und das Maximum wird angenommen für:

$$\arg \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Corr}(Y, \mathbf{a}^T \mathbf{X}) = \frac{\mathbf{V}_{\mathbf{xx}}^{-1/2} \mathbf{V}_{\mathbf{xy}}}{\sqrt{\mathbf{V}_{\mathbf{yx}} \mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{yx}}}}.$$

Also erhalten wir wieder die Regressionsfunktion als Lösung des Maximierungsproblems auf LP_1 .

5.3.2. Vorhersage durch Maximierung der Varianz (PCA-Methode)

In diesem Abschnitt werden wir eine Methode zur Vorhersage vorstellen, die als Hauptkomponentenanalyse bekannt ist, (engl. principal component analysis PCA). Sie ist eines der ältesten Ordinationsverfahren zur Dimensionsreduktion. Sie wurde in der ersten Hälfte des 20. Jahrhunderts im angloamerikanischen Raum entwickelt.

Definition 5.1 (*PCA-Kombination, Koeffizientenvektor, PCA-Vorhersage erster Art*) Puntanen, Styan und Isotalo (2013, Kap.9.4)

- Die Linearkombination $\mathbf{a}_{pca}^T \mathbf{X}$ heißt PCA-Komponente von \mathbf{X} , falls sie die maximale Varianz besitzt.

$$\text{Var}(\mathbf{a}_{pca}^T \mathbf{X}) = \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Var}(\mathbf{a}^T \mathbf{X}). \quad (5.19)$$

- Dabei heißt \mathbf{a}_{pca} der Koeffizientenvektor (Richtungsvektor) der PCA-Komponente
- Die Vorhersage $g_{pca}(\mathbf{X}) = \mu_y + \mathbf{a}_{pca}^T (\mathbf{X} - \boldsymbol{\mu}_x)$ heißt **PCA-Vorhersage erster Art** von Y unter \mathbf{X} in der Klasse LP_1 .

Eigenschaften des PCA-Koeffizientenvektors

(i) **Satz 5.1** *Es gilt*

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Var}(\mathbf{a}^T \mathbf{X}) = \lambda_{\max}(\mathbf{V}_{\mathbf{xx}}) \quad \text{nach Satz A.2,}$$

und die optimale Lösung wird angenommen für $\mathbf{a}_{\text{pca}} = e.v_{\max}(\mathbf{V}_{\mathbf{xx}})$, der zu dem größten Eigenwert λ_{\max} der Matrix $\mathbf{V}_{\mathbf{xx}}$ gehört.

(ii) Mit Hilfe der Projektion können wir die Bestimmung des PCA-Koeffizientenvektors als Minimierungsproblem betrachten, wie der folgende Satz aussagt.

Satz 5.2 *Für $\mathbf{a} \in \mathbb{R}^p$ gilt:*

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Var}(\mathbf{a}^T \mathbf{X}) \iff \min_{\mathbf{a}} \text{Spur}(\text{Cov}[(\mathbf{X} - P_{\mathbf{a}}\mathbf{X})^T(\mathbf{X} - P_{\mathbf{a}}\mathbf{X})])$$

unter der Bedingung $\mathbf{a}^T \mathbf{a} = 1$,

wobei $P_{\mathbf{a}} = \mathbf{a}\mathbf{a}^T$ die Projektionsmatrix ist.

Beweis:

Man kann sich überlegen, dass die Projektion von \mathbf{x} auf \mathbf{a} gleich

$$P_{\mathbf{a}}\mathbf{x} = \frac{\mathbf{a}\mathbf{a}^T}{\mathbf{a}^T \mathbf{a}} \mathbf{x} = \mathbf{a}\mathbf{a}^T \mathbf{x}$$

ist.

$$\begin{aligned} \min_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Spur}(\text{Cov}[(\mathbf{X} - P_{\mathbf{a}}\mathbf{X})^T(\mathbf{X} - P_{\mathbf{a}}\mathbf{X})]) &= \min_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Spur}(\mathbf{V}) - \mathbf{a}^T \mathbf{V}_{\mathbf{xx}} \mathbf{a} \\ &\iff \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Var}(\mathbf{a}^T \mathbf{X}) \end{aligned}$$

(iii) Unter Normalverteilungsannahme mit Erwartungswert gleich 0 gilt

Satz 5.3 *Puntanen, Styan und Isotalo (2013, Kap.9.4)*

Für $\mathbf{a} \in \mathbb{R}^p$ gilt

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Var}(\mathbf{a}^T \mathbf{X}) \iff \min_{\mathbf{a}} \text{MSEP}(\mathbf{X}, E(\mathbf{X}|\mathbf{a}^T \mathbf{X}))$$

unter der Bedingung $\mathbf{a}^T \mathbf{a} = 1$.

Folgerung 5.1 *Für $\mathbf{a} \in \mathbb{R}^p$ gilt*

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Var}(\mathbf{a}^T \mathbf{X}) \iff \min_{\mathbf{a}} \text{MSEP}(\mathbf{X}, \text{BLUP}(\mathbf{X}, \mathbf{a}^T \mathbf{X}))$$

unter der Bedingungen $\mathbf{a}^T \mathbf{a} = 1$.

5.3.3. Vorhersage durch Maximierung der Kovarianz: (PLS-Methode)

Definition 5.2 (*PLS-Kombination, PLS-Koeffizientenvektor, PLS-Vorhersage erster Art*) Liu und Rayens (2007) und Garthwaite (1994)

- Die Linearkombination $\mathbf{a}_{pls}^T \mathbf{X}$ heißt PLS-Komponente von Y unter \mathbf{X} , falls sie die maximale Kovarianz mit Y besitzt.

$$\text{Cov}(Y, \mathbf{a}_{pls}^T \mathbf{X}) = \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Cov}(Y, \mathbf{a}^T \mathbf{X}). \quad (5.20)$$

- Dabei heißt \mathbf{a}_{pls} der Koeffizientenvektor (Richtungsvektor) der PLS-Komponente
- Die Vorhersage $g_{pls}(\mathbf{X}) = \mu_y + \mathbf{a}_{pls}^T (\mathbf{X} - \boldsymbol{\mu}_x)$ heißt **PLS-Vorhersage erster Art** von Y unter \mathbf{X} in der Klasse LP_1 .

Dann gilt der Satz:

Satz 5.4 Die Optimallösung ist gegeben durch:

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Cov}(Y, \mathbf{a}^T \mathbf{X}) = \sqrt{\lambda_{\max}(\mathbf{V}_{xy} \mathbf{V}_{xy}^T)} \quad (\text{nach } A1, A2),$$

und sie wird angenommen für $\mathbf{a}_{pls} = e.v_{\max}(\mathbf{V}_{xy} \mathbf{V}_{xy}^T)$, der zu dem größten Eigenwert der Matrix $\mathbf{V}_{xy} \mathbf{V}_{xy}^T$ gehört.

Bemerkung 5.1 Da $\text{Rang}(\mathbf{V}_{xy} \mathbf{V}_{yx}) = 1$ ist, erhalten wir, dass

$$\mathbf{a}_{pls} = c \mathbf{V}_{xy} \quad \text{mit } c = \frac{1}{\|\mathbf{V}_{xy}\|}. \quad (5.21)$$

Außerdem sind die folgenden Eigenschaften des PLS-Koeffizientenvektors nützlich:

Eigenschaften des PLS-Koeffizientenvektors

- (i) Frank und Friedmann (1993) haben die folgende Beziehung zwischen PLS und PCA-Koeffizient gefunden:

Theorem 5.1 Frank und Friedmann (1993)

$$\begin{aligned} \mathbf{a}_{pls} &= \arg \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \text{Corr}^2(\mathbf{a}^T \mathbf{X}, Y) \text{Var}(\mathbf{a}^T \mathbf{X}) \\ &= \arg \max_{\mathbf{a}^T \mathbf{a} = 1} \text{Cov}^2(\mathbf{a}^T \mathbf{X}, Y) \end{aligned}$$

- (ii) Auf die gleiche Weise können wir beweisen, dass

Theorem 5.2

$$\mathbf{a}_{pls} = \arg \min_{\mathbf{a}^T \mathbf{a} = 1} [Var(Y - \mathbf{a}^T \mathbf{X}) + Spur(Cov(\mathbf{X} - \mathbf{a} \mathbf{a}^T \mathbf{X}))] \quad (5.22)$$

Beweis: Man kann sich überlegen, dass

$$Var(Y - \mathbf{a}^T \mathbf{X}) = Var(Y) + \mathbf{a}^T \mathbf{V}_{xx} \mathbf{a} - 2 Cov(Y, \mathbf{a}^T \mathbf{X}) \quad (5.23)$$

und

$$Spur(Cov(\mathbf{X} - \mathbf{a} \mathbf{a}^T \mathbf{X})) = Spur(Cov[(\mathbf{I} - \mathbf{a} \mathbf{a}^T) \mathbf{X}]) \quad (5.24)$$

$$= Spur(\mathbf{V}_{xx} - \mathbf{V}_{xx} \mathbf{a} \mathbf{a}^T), \text{ da } \mathbf{a}^T \mathbf{a} = 1 \quad (5.25)$$

$$= Spur(\mathbf{V}_{xx}) - \mathbf{a}^T \mathbf{V}_{xx} \mathbf{a} \quad (5.26)$$

Hieraus folgt, dass

$$\arg \min_{\mathbf{a}^T \mathbf{a} = 1} [Var(Y - \mathbf{a}^T \mathbf{X}) + Spur(Cov(\mathbf{X} - \mathbf{a} \mathbf{a}^T \mathbf{X}))] = \arg \max_{\mathbf{a}^T \mathbf{a} = 1} Cov(\mathbf{a}^T \mathbf{X}, Y). \quad (5.27)$$

bzw.

$$\min_{\mathbf{a}^T \mathbf{a} = 1} [Var(Y - \mathbf{a}^T \mathbf{X}) + Spur(Cov(\mathbf{X} - \mathbf{a} \mathbf{a}^T \mathbf{X}))] \Leftrightarrow \max_{\mathbf{a}^T \mathbf{a} = 1} Cov(\mathbf{a}^T \mathbf{X}, Y). \quad (5.28)$$

5.4. Modifizierte partielle kleinste Quadrate

In diesem Abschnitt betrachten wir ein etwas allgemeineres Modell als das, was wir im Abschnitt 5.2 vorgestellt haben, nämlich das Modell 1.3.a mit $q > 1$, d.h.

$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$ eine \mathbb{R}^{p+q} wertige ZV, mit

$\boldsymbol{\mu} \in \mathbb{R}^{p+q}, \mathbf{V} \in PD(p+q)$

$\boldsymbol{\mu}, \mathbf{V}$ bekannt

\mathbf{Y} nicht beobachtbar.

Das Vorgehen der kanonischen Korrelation stellt eine Verallgemeinerung des Multiplen Korrelationskoeffizienten dar. Eine Verallgemeinerung der PLS-Methode wurde von Liu und Rayens (2007) vorgeschlagen. Die Bestimmung des PLS-Koeffizientenvektors in diesem Modell ist gegeben durch das Maximierungsproblem:

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} Cov(\mathbf{a}^T \mathbf{X}, \mathbf{Y}) [Cov(\mathbf{a}^T \mathbf{X}, \mathbf{Y})]^T \text{ bzw.} \quad (5.29)$$

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} \mathbf{a}^T \mathbf{V}_{xy} (\mathbf{a}^T \mathbf{V}_{xy})^T \quad \text{Liu und Rayens (2007)}$$

Im Fall $q = 1$ erhalten wir:

$$\max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} Cov(Y, \mathbf{a}^T \mathbf{X}) \Leftrightarrow \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} [Cov(Y, \mathbf{a}^T \mathbf{X})]^2.$$

Hieraus und basierend auf der Überlegung von Abschnitt 5.1 werden wir eine Verallgemeinerung der partiellen kleinsten Quadrate vorschlagen, die es so in der Literatur noch nicht gibt. Die Idee der modifizierten PLS-Methode liegt darin, dass wir statt des Optimierungsproblems (5.29) betrachten:

$$\max_{\mathbf{a} \in \mathbb{R}^p \setminus \{0\}} \text{Spur}[Cov(BLUP(\mathbf{Y}, \mathbf{a}^T \mathbf{X}), \mathbf{Y})] \quad (5.30)$$

Beachte:

- 1) $BLUP(\mathbf{Y}, \mathbf{a}^T \mathbf{X})$ ist \mathbb{R}^q -wertige Zufallsvariable und $Cov(BLUP(\mathbf{Y}, \mathbf{a}^T \mathbf{X}), \mathbf{Y}) \in PD(q)$.
- 2) Die Lösung von (5.30) im Fall $q = 1$ ist gleich der BLUP

Definition 5.3 (PLS1-Kombination, PLS1-Koeffizientenvektor)

- Die Linearkombination $\mathbf{a}_{pls1}^T \mathbf{X}$ heißt PLS1-Komponente von Y unter \mathbf{X} , falls gilt

$$\mathbf{a}_{pls1} = \arg \max_{\mathbf{a} \in \mathbb{R}^p \setminus \{0\}} \text{Spur}[Cov(BLUP(\mathbf{Y}, \mathbf{a}^T \mathbf{X}), \mathbf{Y})] \quad (5.31)$$

- Dabei heißt \mathbf{a}_{pls1} der Koeffizientenvektor (Richtungsvektor) der PLS1-Komponente

Eigenschaften des PLS1-Koeffizientenvektors

(i) Unter der Annahme, dass $E(\mathbf{X}) = 0$ und $E(\mathbf{Y}) = 0$ sind, gilt:

$$\max_{\mathbf{a} \in \mathbb{R}^p \setminus \{0\}} \text{Spur}[Cov(BLUP(\mathbf{Y}, \mathbf{a}^T \mathbf{X}), \mathbf{Y})] = \min_{\mathbf{a} \in \mathbb{R}^p \setminus \{0\}} MSEP(\mathbf{Y}, BLUP(\mathbf{Y}, \mathbf{a}^T \mathbf{X})).$$

Beweis:

Gemäß der Definition der besten erwartungstreuen Vorhersage erhalten wir:

$$BLUP(\mathbf{Y}, \mathbf{a}^T \mathbf{X}) = \hat{\mathbf{Y}}_{\mathbf{a}} = \frac{\mathbf{V}_{yx} \mathbf{a} \mathbf{a}^T}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \mathbf{X} \quad \text{und} \quad Cov(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}}) = \frac{\mathbf{V}_{yx} \mathbf{a} \mathbf{a}^T \mathbf{V}_{xy}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}}.$$

Der mittlere quadratischen Fehler der Vorhersage im mehrdimensionalen Fall ist

$$\begin{aligned} MSEP(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}}) &= E((\mathbf{Y} - \hat{\mathbf{Y}}_{\mathbf{a}})^T (\mathbf{Y} - \hat{\mathbf{Y}}_{\mathbf{a}})) \\ &= \text{Spur} [Cov(\mathbf{Y} - \hat{\mathbf{Y}}_{\mathbf{a}})] + [E(\mathbf{Y} - \hat{\mathbf{Y}}_{\mathbf{a}})]^T [E(\mathbf{Y} - \hat{\mathbf{Y}}_{\mathbf{a}})] \\ &= \text{Spur}[Cov(\mathbf{Y}) + Cov(\hat{\mathbf{Y}}_{\mathbf{a}}) - Cov(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}}) - Cov(\hat{\mathbf{Y}}_{\mathbf{a}}, \mathbf{Y})] + 0 \\ &= \text{Spur} \left[\mathbf{V}_{yy} - \frac{\mathbf{V}_{yx} \mathbf{a} \mathbf{a}^T \mathbf{V}_{xy}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \right]. \end{aligned}$$

und daher

$$\begin{aligned} \min_{\mathbf{a}} \text{Spur} \left[\mathbf{V}_{yy} - \frac{\mathbf{V}_{yx} \mathbf{a} \mathbf{a}^T \mathbf{V}_{xy}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \right] &\Leftrightarrow \max_{\mathbf{a}} \frac{\mathbf{a}^T \mathbf{V}_{xy} \mathbf{V}_{yx} \mathbf{a}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \\ \min_{\mathbf{a}} MSEP(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}}) &\Leftrightarrow \max_{\mathbf{a}} \text{Spur}[Cov(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}})]. \end{aligned}$$

(ii) **Satz 5.5** *Es gilt*

$$\max_{\mathbf{a} \in \mathbb{R}^p \setminus \{0\}} \text{Spur}[\text{Cov}(\text{BLUP}(\mathbf{Y}, \mathbf{a}^T \mathbf{X}), \mathbf{Y})] = \lambda_{\max}(\mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}} \mathbf{V}_{\mathbf{xy}}^T) \quad (\text{nach Satz A.3}),$$

und die optimale Lösung wird angenommen für $\mathbf{a}_{pls1} = e.v_{\max}(\mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}} \mathbf{V}_{\mathbf{xy}}^T)$, der zu dem größten Eigenwert der Matrix $\mathbf{V}_{\mathbf{xx}}^{-1} \mathbf{V}_{\mathbf{xy}} \mathbf{V}_{\mathbf{xy}}^T$ gehört.

Im sechsten Kapitel findet sich eine Simulationsstudie unter der multivariaten Normalverteilungsannahme mit einer autoregressiven Kovarianzmatrix, um die Güte der Methoden: modifizierte partielle kleinste Quadrate, partielle kleinste Quadrate, Hauptkomponenten und kanonische Korrelation, miteinander zu vergleichen, am Beispiel des Vorhersageintervalls. Die Simulationsstudie zeigt die Vorteile der modifizierten partiellen kleinsten Quadrate gegenüber den bekannten Methoden.

5.5. Algorithmus der PLS-Methode: NIPALS-Algorithmus

Wenn der Prognosefehler $Y - \hat{Y}$ orthogonal zu allen in die Prognosefunktion eingehenden Variablen \mathbf{X} steht, können wir in diesem Sinn die Vorhersage nicht mehr verbessern. Auf der anderen Seite leiten wir folgenden Algorithmus her, um die Vorhersage zu verbessern. Dieser Algorithmus wurde ursprünglich von H. Wold 1975 (siehe Krämer (2007)) eingeführt.

Algorithmus: NIPALS-Algorithmus (siehe Krämer (2007), Li, Udén und v. Rosen (2013))

- **Input:**
 \mathbf{X} eine \mathbb{R}^p -wertige ZV (beobachtbar)
 Y eine \mathbb{R} -wertige ZV (nicht beobachtbar)
 \mathbf{X}, Y seien gemeinsam verteilt.

$$\begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix} \text{ eine } \mathbb{R}^p \times \mathbb{R} \text{ wertige ZV, mit}$$

$$E \begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \text{Cov} \begin{pmatrix} \mathbf{X} \\ Y \end{pmatrix} = \begin{pmatrix} \mathbf{V}_{\mathbf{xx}} & \mathbf{V}_{\mathbf{xy}} \\ \mathbf{V}_{\mathbf{yx}} & V_{yy} \end{pmatrix} \in PD(p+1) \text{ bekannt,}$$

Die PLS-Vorhersage erster Art von Y : $g_{pls}(\mathbf{X}) = \mathbf{a}_{pls}^T \mathbf{X}$.

- **Ziel:** Verbesserung von $g_{pls}(\mathbf{X})$
- $\mathbf{X}_1 = \mathbf{X}$
 For $i = 1, \dots, m < p$ do
 - 1: $\mathbf{a}_i = \arg \max_{\mathbf{a}^T \mathbf{a}=1} \text{Cov}(\mathbf{a}^T \mathbf{X}_i, Y) \text{Cov}(Y, \mathbf{a}^T \mathbf{X}_i)$
 - 2: $T_i = \mathbf{a}_i^T \mathbf{X}_i$ latente Komponente
 - 3: $\mathbf{X}_{i+1} = \mathbf{X}_i - \text{BLUP}(\mathbf{X}_i, T_i)$

end For

• **Output:**

$\mathbf{T} = (\mathbf{T}_1, \dots, \mathbf{T}_m)^T$ eine \mathbb{R}^m -wertige ZV (latente Variable)

$\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_m) \in \mathbb{R}^{p \times m}$ (Gewichtsvektoren nach PLS)

$\hat{Y}_{pls} = BLUP(Y, \mathbf{T})$ (Vorhersage von Y) Helland (1988).

Dieser PLS-Algorithmus läuft nun in drei Phasen ab. In der ersten Phase werden wir eine Vorhersage bestimmen, die als PLS-Vorhersage (erster Art) heißt. In der zweiten Phase werden wir die latente Variable \mathbf{T} herleiten. In der dritten Phase werden wir die beste lineare Vorhersage von Y unter \mathbf{T} bestimmen.

Eigenschaften der Komponenten T_i

1. \mathbf{X}_{i+1} und T_i sind orthogonal, da

$$\begin{aligned} Cov(\mathbf{X}_{i+1}, T_i) &= Cov(\mathbf{X}_i - BLUP(\mathbf{X}_i, T_i), T_i) \\ &= \mathbf{0}. \end{aligned}$$

2. Hieraus folgt auch, dass

$$\begin{aligned} Cov(T_{i+1}, T_i) &= Cov(\mathbf{a}_{i+1}^T \mathbf{X}_{i+1}, T_i) \\ &= \mathbf{a}_{i+1}^T Cov(\mathbf{X}_{i+1}, T_i) \\ &= \mathbf{a}_{i+1}^T Cov(\mathbf{X}_i - BLUP(\mathbf{X}_i, T_i), T_i) \\ &= \mathbf{0} \quad \text{Theorem 3.1.} \end{aligned}$$

Beachte

• Es folgt aus Folgerung 3.4, dass

$$Cov(\mathbf{X}_i - BLUP(\mathbf{X}_i, T_i)) + Cov(BLUP(\mathbf{X}_i, T_i)) = Cov(\mathbf{X}_i)$$

ist. Hieraus ergibt sich $Cov(\mathbf{X}_{i+1}) \leq Cov(\mathbf{X}_i)$

• wenn $g(\mathbf{X}) = \mathbf{a}_{pls}^T \mathbf{X}$ eine BLUP für Y ist, folgt

$$\begin{aligned} \hat{Y}_{pls} &= BLUP(Y, \mathbf{a}_{pls}^T \mathbf{X}) \\ &= BLUP(Y, BLUP(\mathbf{Y}, \mathbf{a}_{pls}^T \mathbf{X})) \\ &= BLUP(Y, \mathbf{X}) \end{aligned}$$

• MSEP von \hat{Y}_{pls} ist gegeben durch:

$$MSEP(Y, \hat{Y}_{pls}) = V_{yy} - \mathbf{V}_{yx} \mathbf{A} (\mathbf{A}^T \mathbf{V}_{xx} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{V}_{xy}$$

5. Partielle kleinste Quadrate

- Im Fall $m=1$ erhalten wir:

(i) MSEP von \hat{Y}_{pls} ist gegeben durch:

$$MSEP(Y, \hat{Y}_{pls}) = V_{yy} - \frac{\mathbf{a}_{pls}^T \mathbf{V}_{xy} \mathbf{V}_{yx} \mathbf{a}_{pls}}{\mathbf{a}_{pls}^T \mathbf{V}_{xx} \mathbf{a}_{pls}}$$

(ii) Der multiple Korrelationskoeffizient von \hat{Y}_{pls} ist gegeben durch:

$$Corr^2(Y, \hat{Y}_{pls}) = \frac{\mathbf{a}_{pls}^T \mathbf{V}_{xy} \mathbf{V}_{yx} \mathbf{a}_{pls}}{V_{yy} \mathbf{a}_{pls}^T \mathbf{V}_{xx} \mathbf{a}_{pls}}$$

5.6. Schlussbemerkung und Interpretation

In diesem Abschnitt betrachten wir das Modell 1.3.b \mathbb{R}^3 -wertige normal-verteilte Zufallsvariable mit autoregressiver Varianz-Kovarianz-Matrix \mathbf{V} :

$$\begin{pmatrix} X_1 \\ X_2 \\ Y \end{pmatrix} \sim N_3 \left(\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho & \rho^2 \\ \rho & 1 & \rho \\ \rho^2 & \rho & 1 \end{pmatrix} \right) \quad \text{mit } \rho \text{ bekannt und } |\rho| < 1.$$

- Der Koeffizientenvektor und **die PLS-Vorhersage erster Art** von Y auf \mathbf{X} sind gegeben durch:

$$\mathbf{a}_{pls} = \frac{1}{\sqrt{1+\rho^2}} \begin{pmatrix} \rho \\ 1 \end{pmatrix}$$

$$g_{pls}(\mathbf{X}) = \frac{1}{\sqrt{1+\rho^2}} (\rho X_1 + X_2) \quad (\text{nach Satz A.1})$$

- Der Koeffizientenvektor und **die PCA-Vorhersage erster Art** von Y auf \mathbf{X} sind gegeben durch:

$$\rho > 0: \quad \mathbf{a}_{pca} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \quad g_{pca}(\mathbf{X}) = \frac{1}{\sqrt{2}} (X_1 + X_2)$$

$$\rho < 0: \quad \mathbf{a}_{pca} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \quad g_{pca}(\mathbf{X}) = \frac{1}{\sqrt{2}} (X_1 - X_2)$$

- Die beste lineare Vorhersage von Y auf \mathbf{X} ist gegeben durch:

$$BLUP(Y, \mathbf{X}) = \hat{Y} = \rho X_2$$

- Die PLS-Vorhersage von Y auf \mathbf{X} ist gegeben durch:

$$\begin{aligned} \hat{Y}_{pls} &= BLUP(Y, \mathbf{a}_{pls}^T \mathbf{X}) \\ &= \frac{\mathbf{V}_{yx} \mathbf{a}_{pls} \mathbf{a}_{pls}^T}{\mathbf{a}_{pls}^T \mathbf{V}_{xx} \mathbf{a}_{pls}} \mathbf{X} \\ &= \frac{1+\rho^2}{3\rho^2+1} (\rho^2 X_1 + \rho X_2) \end{aligned}$$

- Die PCA-Vorhersage von Y auf \mathbf{X} ist gegeben durch:

$$\begin{aligned}\widehat{Y}_{pca} &= BLUP(Y, \mathbf{a}_{pca}^T \mathbf{X}) \\ &= \frac{\mathbf{V}_{yx} \mathbf{a}_{pca} \mathbf{a}_{pca}^T}{\mathbf{a}_{pca}^T \mathbf{V}_{xx} \mathbf{a}_{pca}} \mathbf{X}\end{aligned}$$

und

$$\rho > 0 : \quad \widehat{Y}_{pca} = \frac{1}{2}(\rho X_1 + \rho X_2), \quad \rho < 0 : \quad \widehat{Y}_{pca} = \frac{1}{2}(\rho X_1 - \rho X_2)$$

- Und somit ergibt sich:

$$\begin{aligned}Corr^2(\widehat{Y}, Y) &= \rho^2, \widehat{Y} = BLUP(Y, \mathbf{X}). \\ Corr^2(\widehat{Y}_{pls}, Y) &= \frac{\rho^2(1 + \rho^2)^2}{3\rho^2 + 1} \\ Corr^2(\widehat{Y}_{pca}, Y) &= 0.5 \rho^2(1 + \rho) \quad \rho > 0 \\ Corr^2(\widehat{Y}_{pca}, Y) &= 0.5 \rho^2(1 - \rho) \quad \rho < 0\end{aligned}$$

- Auf die gleiche Weise ergibt sich der mittlere quadratische Fehler:

$$\begin{aligned}MSEP(\widehat{Y}, Y) &= 1 - \rho^2 \\ MSEP(\widehat{Y}_{pls}, Y) &= 1 - \frac{\rho^2(1 + \rho^2)^2}{3\rho^2 + 1} \\ MSEP(\widehat{Y}_{pca}, Y) &= \frac{\mathbf{V}_{yx} \mathbf{a}_{pca} \mathbf{a}_{pca}^T \mathbf{V}_{yx}^T}{\mathbf{a}_{pca}^T \mathbf{V}_{xx} \mathbf{a}_{pca}} \\ MSEP(\widehat{Y}_{pca}, Y) &= 1 - 0.5 \rho^2(1 + \rho), \quad \rho > 0 \\ MSEP(\widehat{Y}_{pca}, Y) &= 1 - 0.5 \rho^2(1 - \rho), \quad \rho < 0.\end{aligned}$$

Insbesondere folgt:

- (a) Ist $|\rho| < 1$, folgt

$$\begin{aligned}MSEP(\widehat{Y}, Y) &\leq MSEP(\widehat{Y}_{pls}, Y) \leq MSEP(\widehat{Y}_{pca}, Y) \\ Corr^2(\widehat{Y}, Y) &\geq Corr^2(\widehat{Y}_{pls}, Y) \geq Corr^2(\widehat{Y}_{pca}, Y)\end{aligned}$$

- (b) Wenn ρ nahe bei 1 liegt, folgt

$$\begin{aligned}MSEP(\widehat{Y}, Y) &= MSEP(\widehat{Y}_{pls}, Y) = MSEP(\widehat{Y}_{pca}, Y) \approx 0 \\ Corr^2(\widehat{Y}, Y) &= Corr^2(\widehat{Y}_{pls}, Y) = Corr^2(\widehat{Y}_{pca}, Y) \approx 1\end{aligned}$$

5. Partielle kleinste Quadrate

(c) Wenn ρ nahe bei -1 liegt, folgt

$$MSEP(\hat{Y}, Y) = MSEP(\hat{Y}_{pls}, Y) = MSEP(\hat{Y}_{pca}, Y) \approx 0$$

$$Corr^2(\hat{Y}, Y) = Corr^2(\hat{Y}_{pls}, Y) = Corr^2(\hat{Y}_{pca}, Y) \approx 1.$$

(d) Im Fall $\rho = 0$ folgt, dass

$$MSEP(\hat{Y}, Y) = MSEP(\hat{Y}_{pls}, Y) = MSEP(\hat{Y}_{pca}, Y) = 1$$

$$Corr^2(\hat{Y}, Y) = Corr^2(\hat{Y}_{pls}, Y) = Corr^2(\hat{Y}_{pca}, Y) = 0.$$

Zur Erklärung ist die konkrete Situation in den Abbildungen (5.1) und (5.2) dargestellt

(5.1) Vergleich des mittleren quadratischen Fehlers in Abhängigkeit von ρ für die Vorhersagen \hat{Y} , \hat{Y}_{pls} und \hat{Y}_{pca} .

(5.2) Vergleich der Güte der Vorhersagen von Y in Abhängigkeit von ρ für die Vorhersagen \hat{Y} , \hat{Y}_{pls} und \hat{Y}_{pca} .

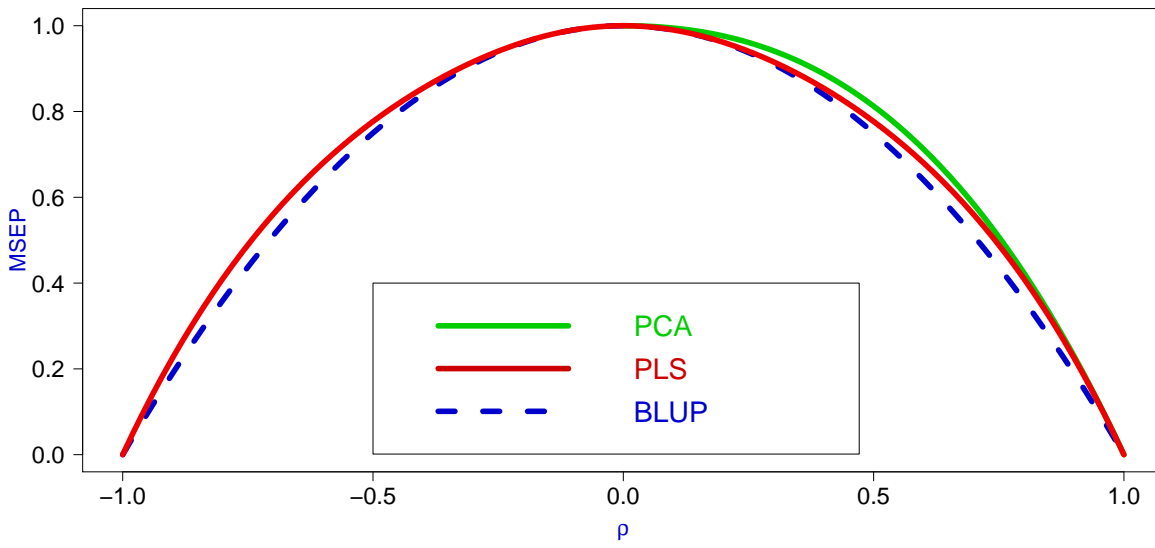


Abbildung 5.1.: Vergleich der MSEP

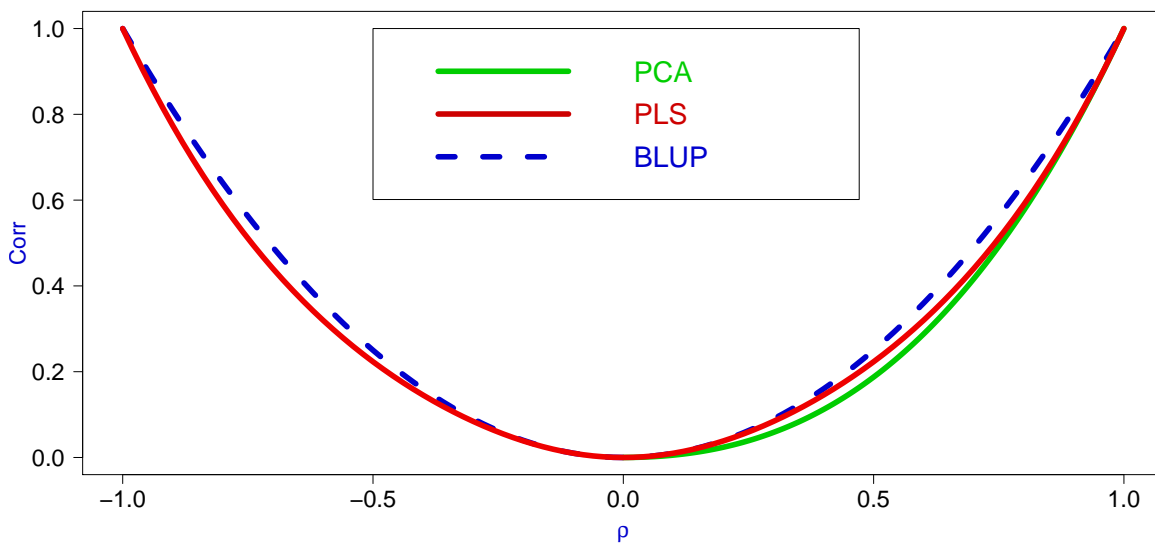


Abbildung 5.2.: Vergleich der Korrelationen

6. Simulationsstudie

In diesem Kapitel werden wir die Güte der Methoden: partielle kleinste Quadrate, Hauptkomponenten und Kanonische Korrelation, modifizierte partielle kleinste Quadrate, am Beispiel der Vorhersageintervalle miteinander durch Simulation vergleichen.

6.1. Modellspezifikation

Wir betrachten wieder das Modell 1.3.b : \mathbb{R}^{p+q} -wertige, normal-verteilte Zufallsvariablen,

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim N_{p+q} \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \right),$$

wobei

- $\mathbf{X} \sim N_p(\mathbf{0}, \mathbf{V}_{xx})$ (beobachtbar)
- \mathbf{Y} eine \mathbb{R}^q -wertige Zufallsvariablen (nicht beobachtbar)
- der Parameter

$$\begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \in PD(p+q),$$

bekannt ist.

6.2. Dimensionsreduktion

Was unter einer statistischen Datenreduktion zu verstehen ist, soll in diesem Abschnitt formalisiert und damit präzisiert werden, d.h. wie man für das Modell

$$\begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim N_{p+q} \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{V}_{xx} & \mathbf{V}_{xy} \\ \mathbf{V}_{yx} & \mathbf{V}_{yy} \end{pmatrix} \right), \quad (6.1)$$

ein Modell

$$\begin{pmatrix} \mathbf{a}^T \mathbf{X} \\ \mathbf{Y} \end{pmatrix} \sim N_{1+q} \left(\begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \mathbf{a}^T \mathbf{V}_{xx} \mathbf{a} & \mathbf{a}^T \mathbf{V}_{xy} \\ \mathbf{V}_{yx} \mathbf{a} & \mathbf{V}_{yy} \end{pmatrix} \right),$$

durch eine Lineartransformation von \mathbf{X} möglichst gut reduziert, mit anderen Worten hat unsere Aufgabe häufig das Ziel, eine niedrigere Dimensionalität zu erreichen, indem eine Lineartransformationen, benutzt wird.

Vorhersage von \mathbf{Y} im reduzierten Modell

$$\begin{aligned}\widehat{\mathbf{Y}}_{\mathbf{a}} &= BLUP(\mathbf{Y}, Z), \text{ mit } Z = \mathbf{a}^T \mathbf{X} \\ &= E(\mathbf{Y}|Z) \\ &= \frac{\mathbf{V}_{\mathbf{y}\mathbf{x}} \mathbf{a} \mathbf{a}^T}{\mathbf{a}^T \mathbf{V}_{\mathbf{x}\mathbf{x}} \mathbf{a}} \mathbf{X} \text{ nach (5.6)}\end{aligned}$$

Eigenschaften von $\widehat{\mathbf{Y}}_{\mathbf{a}}$:

Aus den Verteilungsergebnissen für (mehrdimensionale) normalverteilte Zufallsvariablen erhalten wir, dass

1. $\mathbf{Y} - \widehat{\mathbf{Y}}_{\mathbf{a}} \sim N(\mathbf{0}, \mathbf{C}_{\mathbf{a}})$ das heißt $\mathbf{C}_{\mathbf{a}}^{-1/2}(\mathbf{Y} - \widehat{\mathbf{Y}}_{\mathbf{a}}) \sim N(\mathbf{0}, \mathbf{I}_q)$, wobei

$$\mathbf{C}_{\mathbf{a}} = Cov(\mathbf{Y} - \widehat{\mathbf{Y}}_{\mathbf{a}}) = \mathbf{V}_{\mathbf{y}\mathbf{y}} - \frac{\mathbf{V}_{\mathbf{y}\mathbf{x}} \mathbf{a} \mathbf{a}^T \mathbf{V}_{\mathbf{x}\mathbf{y}}}{\mathbf{a}^T \mathbf{V}_{\mathbf{x}\mathbf{x}} \mathbf{a}} \in PD(q). \quad (6.2)$$

2. Hieraus ergibt sich, dass

$$(\mathbf{Y} - \widehat{\mathbf{Y}}_{\mathbf{a}})^T \mathbf{C}_{\mathbf{a}}^{-1} (\mathbf{Y} - \widehat{\mathbf{Y}}_{\mathbf{a}}) \sim \chi_q^2$$

3. $\frac{Y_i - \widehat{Y}_{\mathbf{a},i}}{\sqrt{C_{\mathbf{a},ii}}} \sim N(0, 1)$, $C_{\mathbf{a},ii} = \mathbf{e}_i^T \mathbf{C}_{\mathbf{a}} \mathbf{e}_i$

- Ein $(1 - \alpha)$ Vorhersageellipsoid für \mathbf{Y} ist gegeben durch:

$$\mathcal{I}_{Sch}(\mathbf{x}) = \{\mathbf{y} \in \mathbb{R}^q : (\mathbf{y} - \widehat{\mathbf{Y}}_{\mathbf{a}}(\mathbf{x}))^T \mathbf{C}_{\mathbf{a}}^{-1} (\mathbf{y} - \widehat{\mathbf{Y}}_{\mathbf{a}}(\mathbf{x})) \leq \chi_{q,1-\alpha}^2\}, \quad (6.3)$$

wobei $\mathcal{I}(\mathbf{x})$ ein Ellipsoid in \mathbb{R}^q mit dem Mittelpunkt $\widehat{\mathbf{Y}}_{\mathbf{a}}(\mathbf{x})$.

Methode zur Bestimmung der DR-Richtung \mathbf{a} **(i) Methode der Hauptkomponente**

\mathbf{a}_{pca} heißt Koeffizientenvektor der Hauptkomponente wenn

$$\mathbf{a}_{pca} = \arg \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} Var(\mathbf{a}^T \mathbf{X})$$

gilt. Die Lösung ist

$$\mathbf{a}_{pca} = e.v_{max}(\mathbf{V}_{\mathbf{x}\mathbf{x}}) \text{ (nach Satz A2)} \quad (6.4)$$

(ii) Methode der partiellen kleinsten Quadrate

\mathbf{a}_{pls} heißt PLS-Koeffizientenvektor wenn

$$\mathbf{a}_{pls} = \arg \max_{\mathbf{a}: \mathbf{a}^T \mathbf{a} = 1} Cov(\mathbf{a}^T \mathbf{X}, \mathbf{Y}) [Cov(\mathbf{a}^T \mathbf{X}, \mathbf{Y})]^T$$

gilt. Die Lösung ist

$$\mathbf{a}_{pls} = e.v_{max}(\mathbf{V}_{\mathbf{x}\mathbf{y}} \mathbf{V}_{\mathbf{y}\mathbf{x}}) \text{ (nach Satz A2)} \quad (6.5)$$

(iii) Methode der modifizierten partiellen kleinsten Quadrate

\mathbf{a}_{pls1} heißt modifizierter PLS-Koeffizientenvektor wenn

$$\mathbf{a}_{pls1} = \arg \max_{\mathbf{a}} \frac{\mathbf{a}^T \mathbf{V}_{xy} \mathbf{V}_{xy}^T \mathbf{a}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}}.$$

Die Lösung ist gegeben durch:

$$\mathbf{a}_{pls1} = e.v_{max}(\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{V}_{yx}) \quad (\text{nach Satz A3}) \quad (6.6)$$

(v) Methode der Kanonischen Korrelation

\mathbf{a}_{cc} heißt CC-Koeffizientenvektor wenn

$$\mathbf{a}_{cc} = \arg \max_{\mathbf{a}} \frac{\mathbf{a}^T \mathbf{V}_{xy} \mathbf{V}_{yy}^{-1} \mathbf{V}_{xy}^T \mathbf{a}}{\mathbf{a}^T \mathbf{V}_{xx} \mathbf{a}} \quad (\text{Satz A.3}).$$

gilt. Die Lösung ist

$$\mathbf{a}_{cc} = e.v_{max}(\mathbf{V}_{xx}^{-1} \mathbf{V}_{xy} \mathbf{V}_{yy}^{-1} \mathbf{V}_{yx}), \quad (6.7)$$

der zum größten Eigenwert gehört.

6.3. Simulationsdesign und Resultate

In diesem Abschnitt werden das Design der Simulationsstudie, die verwendeten Modelle sowie die erhaltenen Resultate zusammenfassend beschrieben.

Das Ziel ist:

- Vergleich des mittleren quadratischen Fehlers zweier Vorhersagen (MSEP) $\widehat{\mathbf{Y}}$ und $\widehat{\mathbf{Y}}_{\mathbf{a}}$ für \mathbf{Y} ,

$$\widehat{\mathbf{Y}} \text{ BUP im vollen Modell} \quad \widehat{\mathbf{Y}}_{\mathbf{a}} \text{ BUP im reduzierten Modell.}$$

- Vergleich zweier Vorhersageintervalle für \mathbf{Y}

Datenerzeugung

$$\mathbf{X} \sim N_2(\mathbf{0}, \mathbf{V}_{xx}) \quad (\text{beobachtbar}) \quad \mathbf{Y} \sim N_2(\mathbf{0}, \mathbf{V}_{yy}) \quad (\text{nicht beobachtbar})$$

(a) Resultate für den Fall: AR(1)-Modell

$$\mathbf{V}_{xx} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}, \mathbf{V}_{yy} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}, \mathbf{V}_{xy} = \begin{pmatrix} \rho^2 & \rho^3 \\ \rho & \rho^2 \end{pmatrix} \quad \text{mit } \rho \text{ bekannt.}$$

- 1) Die beste lineare Vorhersage von Y unter X im vollen Modell und der zugehörige mittlere quadratische Fehler sind:

$$\widehat{Y} = (\rho, \rho^2)^T X_2, \quad MSEP(Y, \widehat{Y}) = 2 - \rho^4 - \rho^2$$

6. Simulationsstudie

- 2) Die beste lineare Vorhersage von Y unter X im reduzierten Modell und der zugehörige mittlere quadratische Fehler sind:

$$\hat{\mathbf{Y}}_{\mathbf{a}} = \frac{\mathbf{V}_{\mathbf{y}\mathbf{x}}\mathbf{a}\mathbf{a}^T}{\mathbf{a}^T\mathbf{V}_{\mathbf{x}\mathbf{x}}\mathbf{a}}\mathbf{X}, \quad MSEP(\mathbf{Y}, \hat{\mathbf{Y}}_{\mathbf{a}}) = \text{Spur}\left(\mathbf{V}_{\mathbf{y}\mathbf{y}} - \frac{\mathbf{V}_{\mathbf{y}\mathbf{x}}\mathbf{a}\mathbf{a}^T\mathbf{V}_{\mathbf{x}\mathbf{y}}}{\mathbf{a}^T\mathbf{V}_{\mathbf{x}\mathbf{x}}\mathbf{a}}\right).$$

Vergleich MSEP der Methoden PCA, PLS, PLS1 und CC:

ρ	BLUP	PCA	PLS	PLS1	CC
.9	0.533900	0.607205	0.599682	0.533900	0.533900
.7	1.269900	1.379415	1.343767	1.269900	1.269900
.3	1.901900	1.936235	1.908226	1.901900	1.901900
.1	1.989900	1.994445	1.989997	1.989900	1.989900

(b) **Resultate für den Fall: Zufällige Kovarianz-Matrizen**

Vergleich MSEP der Methoden PCA, PLS, PLS1 und CC:

BLUP	PCA	PLS	PLS1	CC
4.627538	5.292095	4.737850	4.707023	4.708533
1.622197	2.010749	1.781372	1.624496	1.624543
0.9217632	1.0693525	0.9836368	0.9743091	0.9827330
1.474965	2.592202	1.575064	1.547850	1.547962
2.866695	2.904931	2.867041	2.866729	2.866739
3.818612	4.045535	3.898841	3.826117	3.830563

(c) **Beurteilung der Resultate**

Zur Gütebeurteilung der Methode der modifizierten partiellen kleinsten Quadrate werden die Methode der kanonischen Korrelation und die beste erwartungstreue Vorhersage verwendet. Insgesamt stellt die Methode der modifizierten partiellen kleinsten Quadrate in den betrachteten Situationen eine Verbesserung der PLS-Methode dar.

Die folgenden graphischen Darstellungen veranschaulichen die Resultate der Simulationsstudie des AR(1)-Modell.

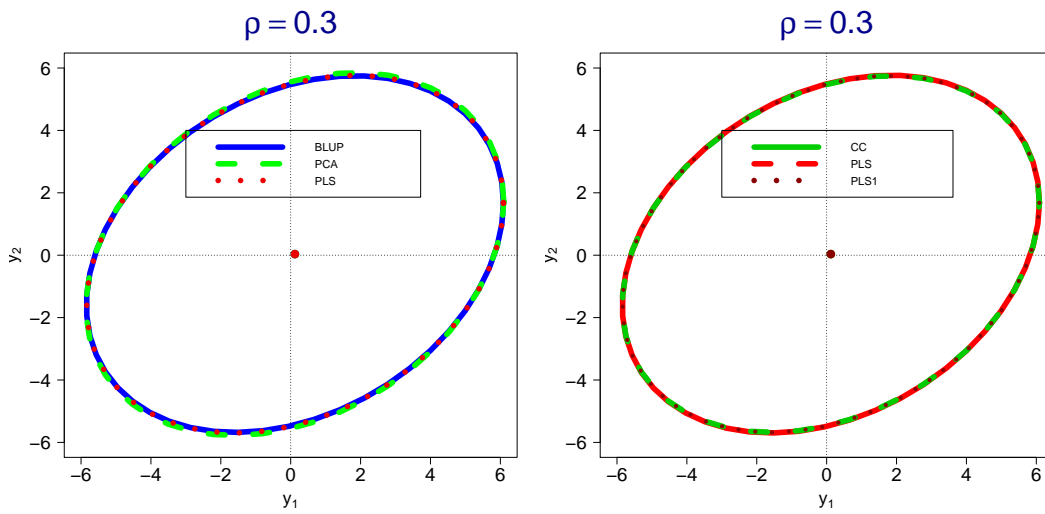


Abbildung 6.1.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a im AR(1)-Modell mit $\rho = .3, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

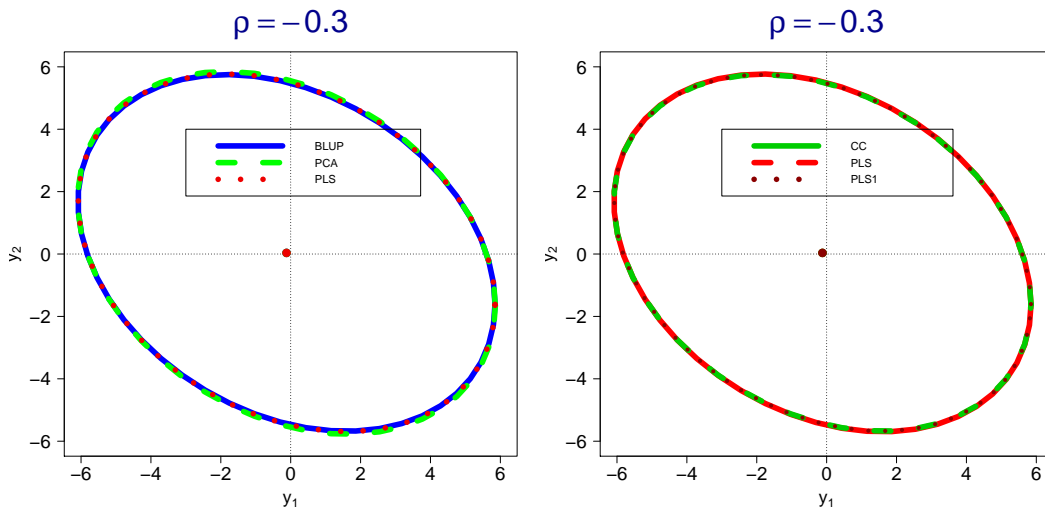


Abbildung 6.2.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a mit $\rho = -0.3, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

6. Simulationsstudie

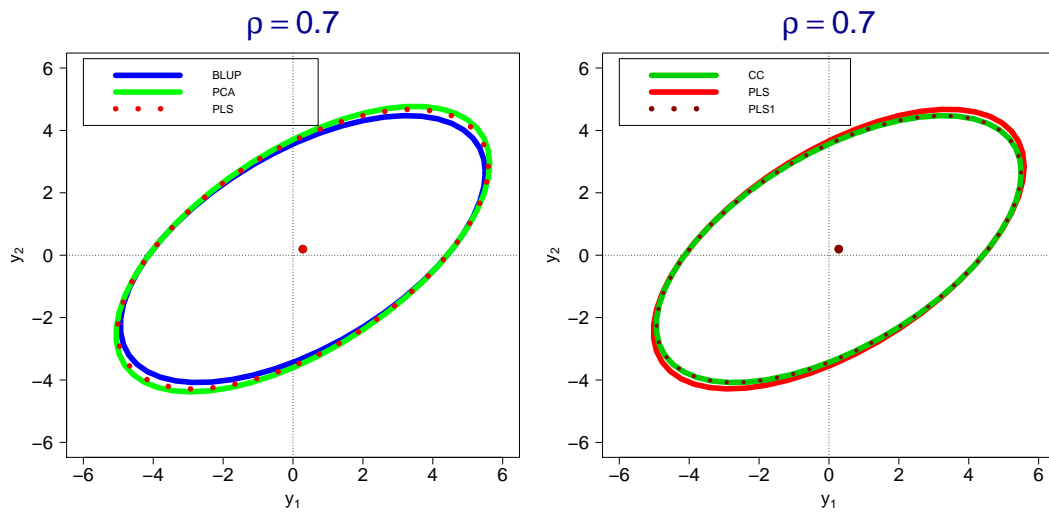


Abbildung 6.3.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a mit $\rho = .7, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

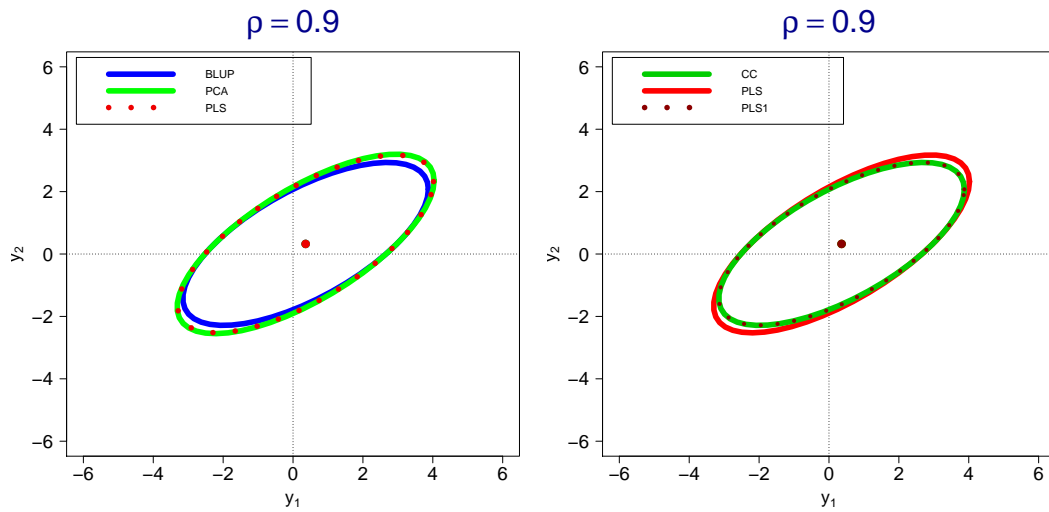


Abbildung 6.4.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a mit $\rho = 0.9, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

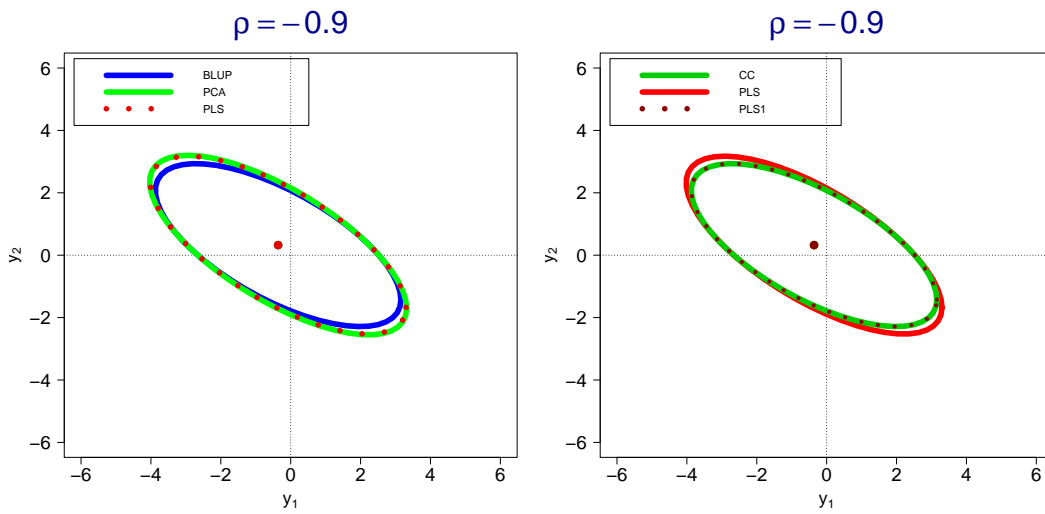


Abbildung 6.5.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a mit $\rho = -0.9, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

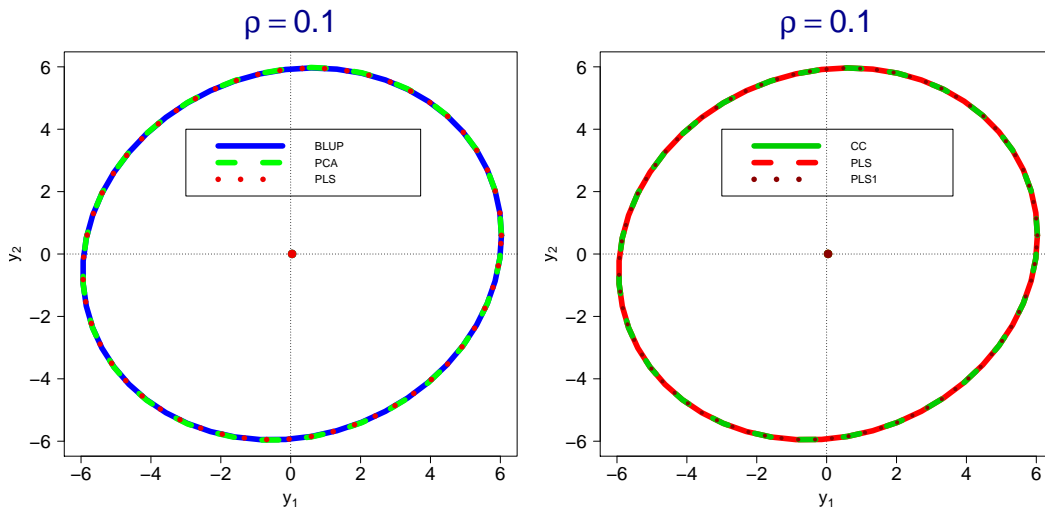


Abbildung 6.6.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a mit $\rho = 0.1, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

6. Simulationsstudie

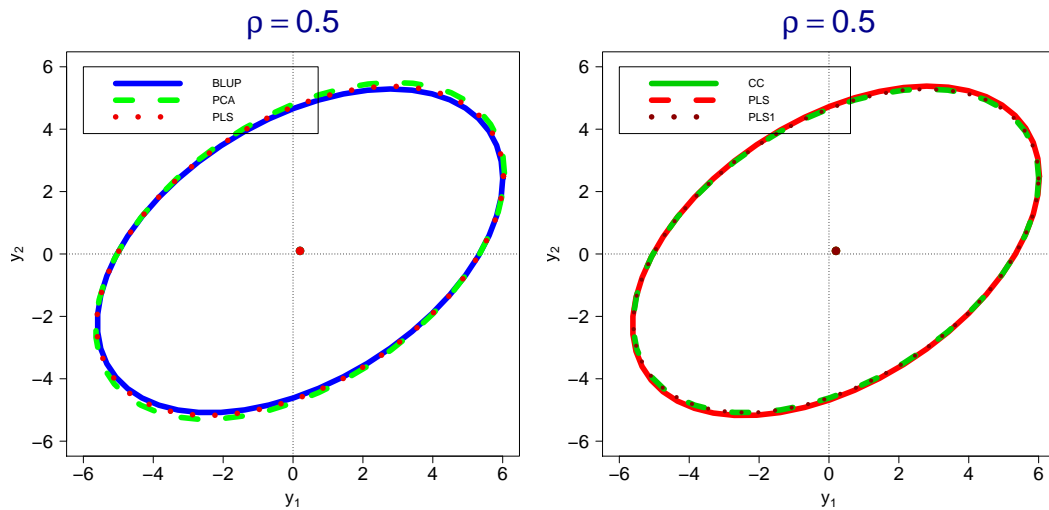


Abbildung 6.7.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a mit $\rho = 0.5, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

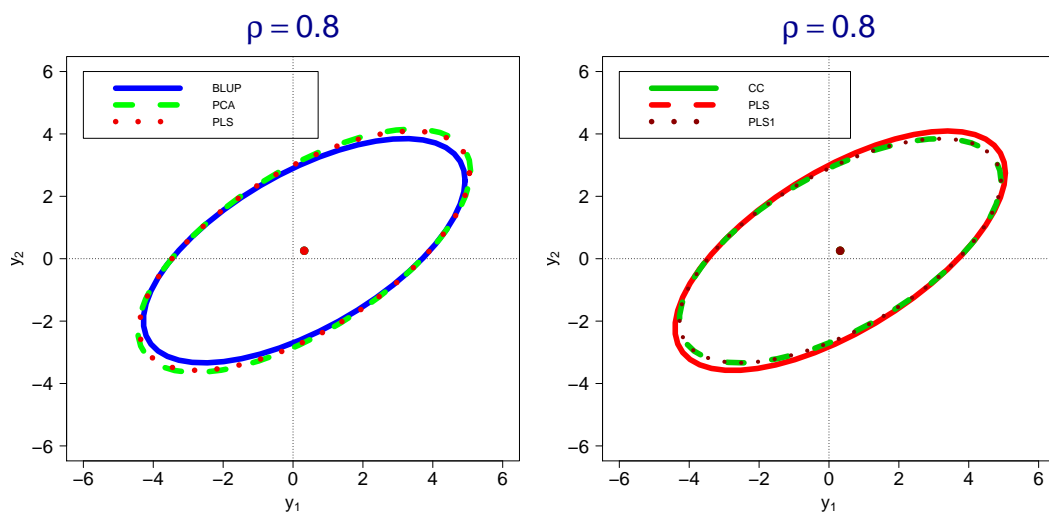


Abbildung 6.8.: Vergleich der Vorhersageellipse für \mathbf{Y} und \mathbf{Y}_a mit $\rho = 0.8, \alpha = .05$. Links(blau: BLUP, grün: PCA, rot: PLS). Rechts(grün: CC, dunkelrot: modifizierte-PLS, rot: PLS)

7. Ausblick

Aufgrund der theoretischen Resultate und der numerischen Vergleiche können wir sagen, dass die PLS-Methode somit gewisse Ähnlichkeiten mit der Hauptkomponentenanalyse (PCA) und der multiplen Korrelation (MK) hat. Bei der PCA-Methode werden ebenfalls zueinander orthogonale $\mathbf{a}^T \mathbf{a} = 1$ neue Variablen als Linearkombination aus den ursprünglichen Variablen bestimmt. Im Unterschied zur PLS werden bei der PCA nicht die Kovarianz (Korrelation), sondern nur die Varianzen maximiert, d.h. bei der PLS wird im Gegensatz zur PCA hier der Zusammenhang mit den Komponenten berücksichtigt. Der mittlere quadratische Fehler der modifizierten PLS ist kleiner als der mittlere quadratische Fehler der normalen PLS. Ebenfalls ist der mittlere quadratische Fehler der modifizierten PLS kleiner als der mittlere quadratische Fehler der CC. Zusammenfassend können wir hinsichtlich der Güte der modifizierten PLS-Methode feststellen, dass die modifizierte PLS-Methode besser als die normale PLS-Methode ist. Ein Nachteil der modifizierten PLS-Methode ist, dass der mittlere quadratische Fehler der modifizierten PLS-Methode groß im Gegensatz zur besten erwartungstreuen Vorhersage sein kann. Ein wesentlicher Vorteil von PLS ist, dass viele Variablen verarbeitet werden können. Es ist sogar möglich mit weniger Versuchen auszukommen, als Variablen vorhanden sind. Insgesamt kann gefolgert werden, dass alle Methoden (BLUP, PLS, PCA, PLS1, CC) gewisse Ähnlichkeiten haben, sodass alle den mittleren quadratischen Fehler der Vorhersage verringern. Natürlich bleiben Ansatzpunkte für weitere Forschung. Eine Aufgabe für die Zukunft besteht zum Beispiel darin, die Theorie der Vorhersage für SUR-Modell (Seemingly Unrelated Regressions) und Growth curve-Modell zu entwickeln. Schließlich kann als Anwendung in der Schätztheorie versucht werden, verbesserte Schätzer im Sinne der Verringerung der Varianz zu finden.

A. Lösungen einiger Maximierungsprobleme

In diesem Anhang geben wir eine kompakte Einführung in Optimierungsprobleme, die auf Matrizen basiert. Für Beweise und weitere Eigenschaften verweisen wir auf die Literatur, z.B. Puntanen, Styan und Isotalo (2013, Kap.22.18,22.25), Fujikoshi, Ulyanov und Shimizu (2010, Anhang. A2).

Satz A.1 Seien \mathbf{B} eine positiv definite symmetrische $n \times n$ -Matrix und $\mathbf{x} \in \mathbb{R}^n$

$$(a) \quad \max_{\mathbf{a}^T \mathbf{B} \mathbf{a} = 1} \mathbf{a}^T \mathbf{x} = \sqrt{\mathbf{x}^T \mathbf{B}^{-1} \mathbf{x}}$$

$$(b) \quad \arg \max_{\mathbf{a}^T \mathbf{B} \mathbf{a} = 1} \mathbf{a}^T \mathbf{x} = \frac{\mathbf{B}^{-1} \mathbf{x}}{\sqrt{\mathbf{x}^T \mathbf{B}^{-1} \mathbf{x}}}.$$

Satz A.2 Sei \mathbf{B} eine symmetrische $n \times n$ -Matrix, dann folgt

$$\max_{\mathbf{a} \in \mathbb{R}^n \setminus \{0\}} \frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{a}} = \lambda_{\max}(\mathbf{B}) \quad \text{und} \quad \arg \max_{\mathbf{a} \in \mathbb{R}^n \setminus \{0\}} \frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{a}} = e.v_{\max}(\mathbf{B})$$

Satz A.3 (Extrema quadratischer Formen unter quadratischer Nebenbedingung)

Seien \mathbf{A} eine symmetrische $n \times n$ -Matrix und \mathbf{B} eine positiv definite $n \times n$ -Matrix. Sei λ_n der größte Eigenwert von $\mathbf{B}^{-1} \mathbf{A}$. Dann gilt

$$\max_{\mathbf{a} \in \mathbb{R}^n \setminus \{0\}} \frac{\mathbf{a}^T \mathbf{A} \mathbf{a}}{\mathbf{a}^T \mathbf{B} \mathbf{a}} = \lambda_n \quad \text{und} \quad \max_{\mathbf{a}: \mathbf{a}^T \mathbf{B} \mathbf{a} = 1} \mathbf{a}^T \mathbf{A} \mathbf{a} = \lambda_n,$$

und die Maximums werden jeweils in einem Eigenvektor der Matrix $\mathbf{B}^{-1} \mathbf{A}$ zu dem Eigenwert λ_n erreicht.

Auf ähnliche Weise, wie im Satz A.1 erhalten wir

Satz A.4 Seien \mathbf{B} eine positiv definite $n \times n$ -Matrix und $\mathbf{x} \in \mathbb{R}^n$. Dann ist

$$(a) \quad \max_{\mathbf{a} \in \mathbb{R}^n \setminus \{0\}} \frac{(\mathbf{a}^T \mathbf{x})^2}{\mathbf{a}^T \mathbf{B} \mathbf{a}} = \mathbf{x}^T \mathbf{B}^{-1} \mathbf{x}$$

$$(b) \quad \arg \max_{\mathbf{a} \in \mathbb{R}^n \setminus \{0\}} \frac{(\mathbf{a}^T \mathbf{x})^2}{\mathbf{a}^T \mathbf{B} \mathbf{a}} = \frac{\mathbf{B}^{-1} \mathbf{x}}{\mathbf{x}^T \mathbf{B}^{-1} \mathbf{x}}.$$

Literaturverzeichnis

- [1] Alvin, R.C. (2008): Linear models in statistics. Wiley New York.
- [2] Christensen, R. (2011): Plane answers to complex questions, Fourth edition. Springer New York.
- [3] Christensen, R., Lin, Y. (2013): Linear models that allow perfect estimation, *stat Papers* 54:695-708.
- [4] Frank, I., Friedman J. A. (1993): Statistical view of some chemometrics regression tools, *Technometrics* 35:109-148.
- [5] Cook, R. D. (1994): On the Interpretation of Regression Plots, *journal of the American Statistical Association*, 89:177-189.
- [6] Fujikoshi, Y., Ulyanov, V. V., Shimizu, R. (2010): *Multivariate Statistics: High Dimensional and Large-Sample Approximations*. Wiley New York.
- [7] Garthwaite, P. H. (1994): An interpretation of partial least squares. *Journal of American Statistical Association*, 89:122-127.
- [8] Grüning, M. (2005): *Untersuchungen zur Diskriminanzanalyse mit hochdimensionalen Daten*, Magdeburg, Univ., Fak. für Mathematik, Diss.
- [9] Helland, I. S. (1988): On the structure of partial least squares regression. *Comm. Statist. Simulation Comput.* 17(2):581-607.
- [10] Helland, I. S. (1990): Partial Least Squares Regression and Statistical Models. *Scand. J. Statist.* 17:97- 114.
- [11] Hocking, R. R. (2005): *Methods and Applications of Linear Models : Regression and the Analysis of Variance*, second edition. Wiley New York.
- [12] Isotalo, J., Puntanen, S. (2009): A note on the equality of the OLSE and the BLUE of the parametric function in the general Gauss-Markov model. *Stat Papers* 50:185:193.
- [13] Krämer, N. (2007): An overview on the shrinkage properties of partial least squares regression. *Journal of Computational.* 22(2): 249-273
- [14] Kreiß, J. P., Neuhaus, G. (2006): *Einführung in die Zeitreihenanalyse*, Springer Berlin.

LITERATURVERZEICHNIS

- [15] Li, K-C. (1991): Siliced Inverse Regression for Dimension Reduction, journal of the American Statistical Association, 86:316-342.
- [16] Li, Y., Udén, P., v. Rosen, D. (2013) : A two-step PLS inspired method for linear prediction with group effect. Sankhya A. 75(1): 96-117
- [17] Liu, Y., Rayens, W., (2007): PLS and dimension reduction for classification. Journal of Computational. 22(2): 189-208
- [18] Puntanen, S., Styan G. P.H., Isotalo, J. (2011): Matrix Tricks for Linear Statistical Models. Springer-Verlag Berlin.
- [19] Puntanen, S., Styan G. P. H., Isotalo, J. (2013): Formulas Useful for Linear Regression Analysis and Related Matrix Theory. Springer-Verlag Berlin.
- [20] Rao, C. R., Shalabh, Toutenburg, H., Heumann. C. (2008): Linear Models and Generalizations : least squares and alternatives, third edition. Springer Berlin.
- [21] Rao, C. R. (1973): Linear Statistical Inference and its Applications, second edition. Wiley New York.
- [22] Schlittgen, R. (2009): Multivariate Statistik. Oldenbourg München
- [23] Schmidt, D. K. (1988): Prediction in the linear model: A direct approach. Metrika. 48(2): 141-147
- [24] Seber, G. A. (1984): Multivariate Observations. Wiley New York.
- [25] Sengupta, D. S., Jammalamadaka, R. S. (2003): Linear Models An Integrated Approach. World Scientific Publishing Company London.
- [26] Shao, J. (2003): Mathematical statistics, second edition. Springer New York.
- [27] Witting, H. (1985): Mathematische Statistik I. Teubner Stuttgart.
- [28] Wold, H. (1966): Estimation of principal components and related models by iterative least squares. In P.R. Krishnaiah (ed.) Multivariate Analysis, 391-420. New York: Academic Press.