

Multivariate Statistik

Vorlesungsskriptum

Univ.-Prof. Dipl.-Ing. Dr.techn.
Peter Filzmoser

Institut für Stochastik und Wirtschaftsmathematik

Wien, im Oktober 2018

Die Vervielfältigung des Skriptums oder von Teilen des Skriptums ist nur mit Genehmigung des Autors zulässig.

Vorwort

Unter *multivariater Analyse* versteht man die Untersuchung mehrerer Variablen, die an mehreren Objekten simultan beobachtet wurden. Das notwendige Instrumentarium dafür nennt man *multivariate Methoden*. *Multivariate Statistik* ist somit die Erweiterung der *univariaten* (eindimensionalen) Statistik auf mehrere Dimensionen.

Ein Beispiel von Daten, mit denen eine multivariate Analyse durchgeführt werden kann, sind Prüfungsergebnisse von Studenten in mehreren Gegenständen. Die Objekte oder Beobachtungen sind die Studenten; ihre Leistungen wurden jeweils für die Gegenstände (Variablen) ermittelt. Der gravierende Unterschied zur *univariaten Analyse* ist, dass nicht jede Variable (Prüfungsfach) gesondert untersucht wird, sondern simultan alle Messungen ausgewertet werden. Für unser Beispiel möchte man etwa Aussagen darüber machen, ob die Unterrichtsgegenstände in Gruppen eingeteilt werden können. Die Interpretation der Gruppen könnte dann lauten, dass die Gegenstände, die eine Gruppe bilden, mathematisches Verständnis voraussetzen, während eine andere Gruppe Gegenstände mit Schwerpunkt auf geometrischem Verständnis beinhaltet. Eine interessante Fragestellung wäre auch, in welchen Unterrichtsfächern die Talente eines Studenten liegen und ob gleich talentierte Studenten die selbe schulische Vorbildung haben. Abhängig von der Antwort könnten Rückschlüsse auf verschiedene Typen von Mittelschulen gemacht werden.

Ein umfangreicheres Beispiel wäre, die Gesellschaftsstruktur der europäischen Länder zu ermitteln. Anhand von Wirtschafts-, Sozial-, Umwelt-, Bevölkerungsdaten etc. könnten interessante Fragestellungen mittels multivariater Analysen (hoffentlich) umfassend beantwortet werden. Die Antworten und Rückschlüsse könnten die Grundlage für politische Entscheidungen bilden.

Schon durch diese beiden Beispiele ist erkennbar, dass Fragestellungen, die mehrere Themen oder Gebiete in Betracht ziehen, oft erst mit Hilfe von multivariaten Analysen beantwortet werden können. Notwendige Voraussetzung für sinnvolle Aussagen sowie für eine sinnvolle Anwendung der Methoden ist natürlich eine entsprechende Erhebung der Daten.

Allgemein könnte man sagen, dass die multivariate Analyse folgende Zielsetzungen hat:

- Datenreduktion und Vereinfachung der zu untersuchenden Struktur
- sortieren und gruppieren
- Untersuchung von Abhängigkeiten der Variablen
- Vorhersage
- Testen von Hypothesen

Die Entwicklung der Theorie der multivariaten Analyse begann etwa 1930, wobei es bei den Methoden vorerst noch starke Einschränkungen gab (z.B. multivariate Normalverteilung). Mit der Erfindung des Computers wurden multivariate Verfahren immer attraktiver. Die schnellen Rechner in der heutigen Zeit erlauben es, die Modelle der multivariaten Methoden immer besser der Wirklichkeit anzupassen.

Inhaltsverzeichnis

Vorwort	i
1 Einführung	1
1.1 Einfache grafische Techniken	1
1.2 Notation	6
1.3 Beschreibende Statistik	7
1.4 Eigenwerte und Eigenvektoren	8
1.5 Erwartungswert und Kovarianzmatrix	10
1.6 Die multivariate Normalverteilung	11
1.7 Transformation zu Normalverteilung	13
1.8 Tests, Konfidenzbereiche	15
2 Clusteranalyse	18
2.1 Einleitung	18
2.2 Klassifikationstyp	18
2.3 Bewertungskriterien für Klassifikationen	19
2.3.1 Maße für die Homogenität einer Klasse	19
2.3.2 Maße für die Heterogenität zwischen den Klassen	20
2.3.3 Maße für die Güte einer Klassifikation	21
2.4 Konstruktionsverfahren	22
2.4.1 Partitionen (Partitionierungsmethoden)	22
2.4.2 Hierarchie	23
2.5 Fuzzy Clusterung	25
3 Multivariate lineare Regression	30
3.1 Einführung	30
3.2 Lineare multiple Regression	30
3.3 Der Kleinste-Quadratsummen-Schätzer	31
3.4 Multivariate lineare Regression	35
4 Auszug aus Robuster Statistik	38
4.1 Einleitung	38
4.2 Robuste lineare Regression – Teil 1: Methoden	38
4.2.1 LS-Regression	38
4.2.2 L_1 -Regression	40
4.2.3 LMS-Regression und LTS-Regression	41

4.3	Robuste Schätzung von multivariater Lokation und Kovarianz	46
4.3.1	Allgemeines	46
4.3.2	MVE und MCD Schätzer	46
4.4	Robuste lineare Regression – Teil 2: Diagnostik	50
4.4.1	Hat-Matrix zur Identifikation von Hebelpunkten	50
4.4.2	Robuste Distanz zur Identifikation von Hebelpunkten	51
4.4.3	Diagnostik bei Regression	52
5	Hauptkomponentenanalyse	57
5.1	Einleitung	57
5.2	Bestimmung von Hauptkomponenten in der Population	57
5.3	Geometrische Interpretation von Hauptkomponenten	60
5.4	Hauptkomponenten von Stichproben	61
5.5	Anzahl der relevanten Hauptkomponenten	64
5.5.1	„Faustregeln“	65
5.6	Singulärwertzerlegung	67
5.7	Biplots	68
5.8	Diagnostik	71
6	Faktorenanalyse	75
6.1	Einleitung	75
6.2	Das Faktorenmodell	75
6.2.1	Definition	75
6.2.2	Nichteindeutigkeit der Faktorenladungen	77
6.2.3	Parameterschätzung	77
6.3	Vorgangsweise	78
6.3.1	Hauptfaktorenanalyse	78
6.3.2	Maximum-Likelihood-Methode	82
6.4	Faktorenrotation	85
6.4.1	Orthogonale Rotationsverfahren	87
6.4.2	Schiefwinkelige Rotationsverfahren	88
6.5	Schätzung von Faktorenwerten	90
6.5.1	Gewichtete Kleinste-Quadratsummen-Schätzung	91
6.5.2	Regressionsmethode	91
7	Korrelationsanalyse	95
7.1	Multiple Korrelation	95
7.2	Kanonische Korrelation	97
8	Diskriminanzanalyse	102
8.1	Einleitung	102
8.2	Überlegungen zu Klassifikationsregeln	102
8.3	Der Zweigruppenfall	104
8.3.1	Spezialfall $\Sigma_1 = \Sigma_2 = \Sigma$	105
8.3.2	Spezialfall $\Sigma_1 \neq \Sigma_2$	109
8.3.3	Auswertung der Klassifikation	110

8.3.4	Die Diskriminanzfunktion von Fisher	111
8.4	Klassifikation mehrerer Populationen	114
8.4.1	Die Methode zur Minimierung der EKM	114
8.4.2	Klassifikation bei Normalverteilung	115
8.4.3	Diskriminanzanalyse nach Fisher für den Mehrgruppenfall . .	119
9	Projection Pursuit	123
9.1	Einleitung	123
9.1.1	Hintergrund	123
9.1.2	Definitionen und Notation	123
9.2	Der Projektionsindex von Friedman	124
9.2.1	Einführung	124
9.2.2	Ein- und zweidimensionale Projektion	125
9.2.3	Entfernung einer Struktur	129
9.2.4	Robustheit	131
9.3	Andere Projektionsindices	131
9.3.1	Der Projektionsindex von Friedman und Tukey	131
9.3.2	Der Entropieindex	132
9.3.3	Der Momentindex	132

Kapitel 1

Einführung

1.1 Einfache grafische Techniken

Ein sehr einfaches aber hilfreiches Mittel in der Datenanalyse sind Grafiken. Gerade wegen ihrer Einfachheit werden diese Hilfsmittel oft als „minderwertig“ oder „wenig wissenschaftlich“ eingestuft. Die Erfahrung zeigt aber, dass eine grafische Darstellung viel Einsicht vermittelt über die Qualität und Struktur der Daten.

Eine sehr bekannte Darstellung von Daten im zweidimensionalen Raum ist das *Streudiagramm*. Dabei werden zwei Variablen einander gegenübergestellt und die Beobachtungen in dieses Koordinatensystem eingetragen.

Beispiel 1.1.1 *In einem Artikel des Magazins “Forbes” vom 30. April 1990 wurden von den größten Verlagen Finanzdaten publiziert. Wir betrachten dabei die Variable x , die die Anzahl der Beschäftigten angibt, und die Variable y mit dem Profit pro Beschäftigtem. Abbildung 1.1 zeigt ein Streudiagramm dieser Daten. Es ist allerdings auffällig, dass “Dun & Bradstreet”, die Firma mit den meisten Beschäftigten, eigentlich typisches Verhalten bei Profit pro Beschäftigtem zeigt und somit vom Hauptteil der Daten abweicht. “Time Warner” hat zwar eine typische Anzahl von Beschäftigten, hat aber im Vergleich geringen Profit (sogar negativ) pro Beschäftigtem.*

Der empirische Korrelationskoeffizient zwischen x und y beträgt -0.39 . Wird “Dun & Bradstreet” von der Berechnung weggelassen, erhält man einen Koeffizienten von -0.56 . Nimmt man zusätzlich “Time Warner” heraus, wird der Korrelationskoeffizient -0.50 . Das Beispiel zeigt also gut, dass Beobachtungen, die von der Hauptstruktur der Daten abweichen, relativ starken Einfluss auf statistische Schätzer haben können. Visuell könnte man diese atypischen Beobachtungen gut durch eine zweidimensionale Erweiterung des Boxplots, den sogenannten “Bagplot” erkennen (Rousseeuw et al., 1999). In Abbildung 1.2 werden die Daten mit Hilfe des Bagplots dargestellt. Analog zum univariaten Boxplot werden (robustes) Mittel, der innere 50%-Bereich der Daten (dunkel schattiert), und ein Bereich (hell schattiert), der die Ausreißer von den regulären Beobachtungen abgrenzt, dargestellt. Man erkennt sehr gut die schon vorher vermuteten Beobachtungen als Ausreißer. Durch die Form des dunkel schattierten “bags” bietet diese Darstellung auch ein Maß für die Korrelation.

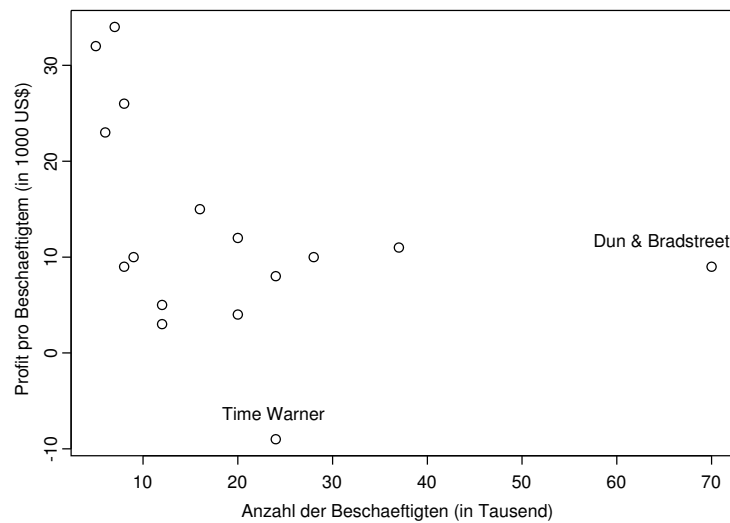


Abbildung 1.1: Streudiagramm für die 2-dimensionalen Finanzdaten

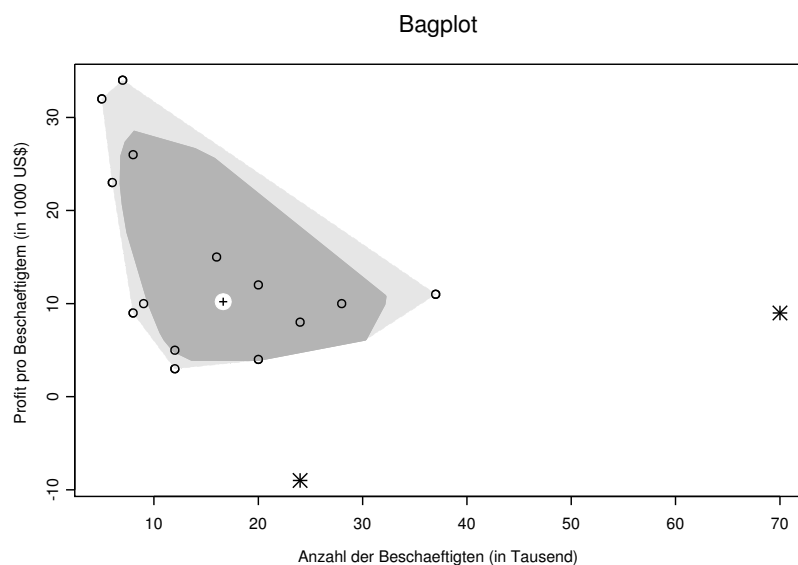


Abbildung 1.2: 2-dimensionale Boxplot-Darstellung der Finanzdaten

Das Streudiagramm kann natürlich auch für mehr als 2-dimensionale Daten eingesetzt werden, wie folgendes Beispiel zeigt:

Beispiel 1.1.2 *Tabelle 1.1 zeigt Daten der Qualität von Papier. Nachdem die Fasern im Papier in einer bestimmten Richtung ausgerichtet sind, werden verschiedene Reißfestigkeiten beobachtet, wenn entlang der Richtung, in der das Papier produziert wird, oder quer dazu gemessen wird. Variable x_1 gibt die Dichte des Papiers in g/cm^3 , Variable x_2 bzw. x_3 geben die Reißfestigkeit entlang der Produktion von der Maschine bzw. quer dazu in Pfund an.*

Diese 3-dimensionalen Daten können wiederum durch Streudiagramme visua-

Tabelle 1.1: Papierqualität; x_1 = Dichte, x_2 = Reißfestigkeit längs, x_3 = Reißfestigkeit quer.

x_1	x_2	x_3	x_1	x_2	x_3	x_1	x_2	x_3
0.801	121.41	70.42	0.832	117.51	71.62	0.822	130.50	80.33
0.824	127.70	72.47	0.796	109.81	53.10	0.822	127.90	75.68
0.841	129.20	78.20	0.759	109.10	50.85	0.843	123.90	78.54
0.816	131.80	74.89	0.770	115.10	51.68	0.824	124.10	71.91
0.840	135.10	71.21	0.759	118.31	50.60	0.788	120.80	68.22
0.842	131.50	78.39	0.772	112.60	53.51	0.782	107.40	54.42
0.820	126.70	69.02	0.806	116.20	56.53	0.795	120.70	70.41
0.802	115.10	73.10	0.803	118.00	70.70	0.805	121.91	73.68
0.828	130.80	79.28	0.845	131.00	74.35	0.836	122.31	74.93
0.819	124.60	76.48	0.822	125.70	68.29	0.788	110.60	53.52
0.826	118.31	70.25	0.971	126.10	72.10	0.772	103.51	48.93
0.802	114.20	72.88	0.816	125.80	70.64	0.776	110.71	53.67
0.810	120.30	68.23	0.836	125.50	76.33	0.758	113.80	52.42
0.802	115.70	68.12	0.815	127.80	76.75			

lisiert werden, indem alle Paare von Variablen als Koordinatenachsen genommen werden (siehe Abbildung 1.3).

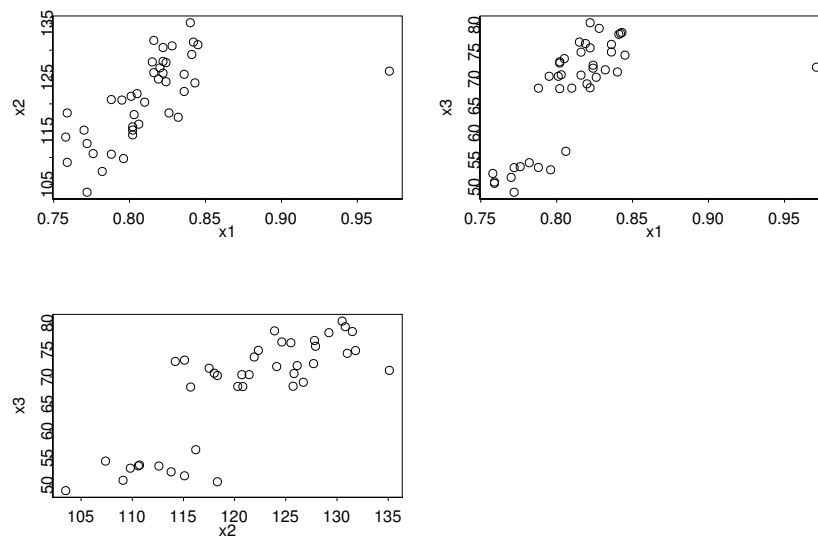


Abbildung 1.3: Streudiagramme der Papierdaten

Die modernere Computersoftware macht es möglich, Daten auch in einem 3-dimensionalen Streudiagramm darzustellen. Durch kontinuierliche Rotation der 3 Koordinatenachsen erhält man visuell einen Eindruck der Datenstruktur. Es ist auch denkbar, mehr als 3 Dimensionen durch ein Streudiagramm darzustellen. Ermöglicht

Tabelle 1.2: Festigkeitsmessung von Platten; x_1 : Messung durch Druckwelle, x_2 : Messung durch Vibration, x_3, x_4 : durch statische Tests.

x_1	x_2	x_3	x_4	x_1	x_2	x_3	x_4
1889	1651	1561	1778	1954	2149	1180	1281
2403	2048	2087	2197	1325	1170	1002	1176
2119	1700	1815	2222	1419	1371	1252	1308
1645	1627	1110	1533	1828	1634	1602	1755
1976	1916	1614	1883	1725	1594	1313	1646
1712	1712	1439	1546	2276	2189	1547	2111
1943	1685	1271	1671	1899	1614	1422	1477
2104	1820	1717	1874	1633	1513	1290	1516
2983	2794	2412	2581	2061	1867	1646	2037
1745	1600	1384	1508	1856	1493	1356	1533
1710	1591	1518	1667	1727	1412	1238	1469
2046	1907	1627	1898	2168	1896	1701	1834
1840	1841	1595	1741	1655	1675	1414	1597
1867	1685	1493	1678	2326	2301	2065	2234
1859	1649	1389	1714	1490	1382	1214	1284

wird dies durch die *Grand-Tour* (Asimov, 1985), einer Methode, die 2-dimensionale Projektionen des p -dimensionalen Raumes berechnet, wobei man durch die p Dimensionen kontinuierlich „hindurchgleitet“. Ein Softwarepaket mit einer Implementierung der Grand-Tour ist *GGobi* (<http://www.ggobi.org>).

Beispiel 1.1.3 Tabelle 1.2 zeigt 4 Festigkeitsmessungen von 30 Platten. Die erste Messung wurde erhalten, indem eine Druckwelle auf die Platte ausgesandt wurde, bei der zweiten Messung wurde die Platte einer Vibration ausgesetzt, und die restlichen beiden Messungen wurden durch statische Tests erhalten.

Abbildung 1.4 zeigt verschiedene Perspektiven der Daten von Tabelle 1.2 in einem 3-dimensionalen Streudiagramm. In Abbildung 1.4a sind die Ausreißer gut erkennbar, Abbildung 1.4b hingegen maskiert die Ausreißer.

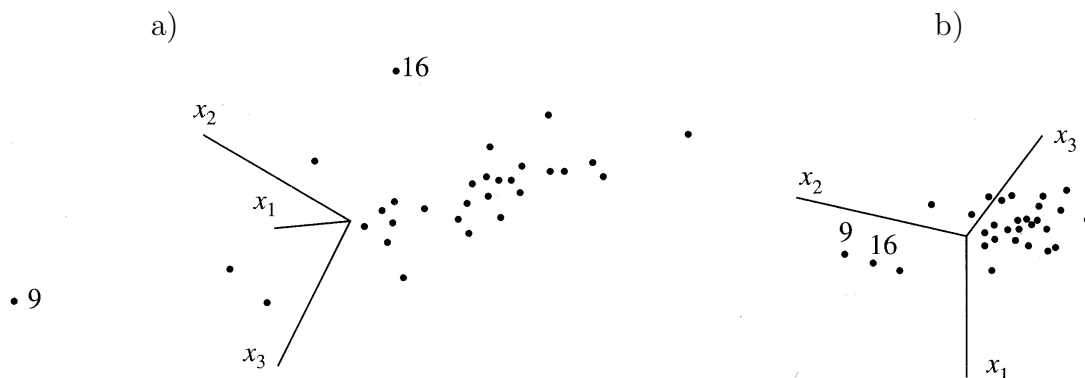


Abbildung 1.4: 3-dimensionale Streudiagramme mit verschiedenen Perspektiven von den Plattendaten.

Selbst für diese einfachen Darstellungen von Daten kann man mit guter Software viel Information extrahieren. Z.B. sind in Abbildung 1.3 deutlich zwei Gruppen und ein Ausreißer erkennbar. Wenn man diese Gruppen markiert, ev. mit verschiedenen Farben, kann man durch Verknüpfungen zu anderen Grafiken (*linking*) diese Gruppen deutlich machen. Besonders interessant ist diese Methode bei Regionaldaten: Die verschiedenfärbigen Gruppen können in einer Landkarte visualisiert werden. Man könnte außerdem Objekte aus anderen Gruppen wegnehmen, und mit der verbleibenden homogenen Gruppe neue Statistiken berechnen.

Es gibt neben den Streudiagrammen noch die *multivariaten Grafiken*. Dabei wird versucht, mehrdimensionale Daten durch bestimmte grafische Symbole zu visualisieren. Bei den *Star-Plots* werden die Variablen als sternförmige Achsen dargestellt, und die Werte für jede Variable eines Objektes in diesem sternförmigen Koordinatensystem eingetragen. Danach werden die eingetragenen Werte benachbarter Variablen verbunden. Es entsteht somit für jedes Objekt ein „sternförmiges Gebilde“, die für relativ gleichartige Objekte sehr ähnlich aussehen. Somit können Gruppen in den Daten sichtbar gemacht werden.

Abbildung 1.5 zeigt einen Star-Plot der Wahlergebnisse der republikanischen Partei von 1856-1976 (31 Werte) für 9 der nordöstlichen Bundesstaaten der USA. Die Werte werden, von rechts beginnend, im Uhrzeigersinn aufgetragen. Man erkennt zum Teil sehr ähnliche Strukturen.

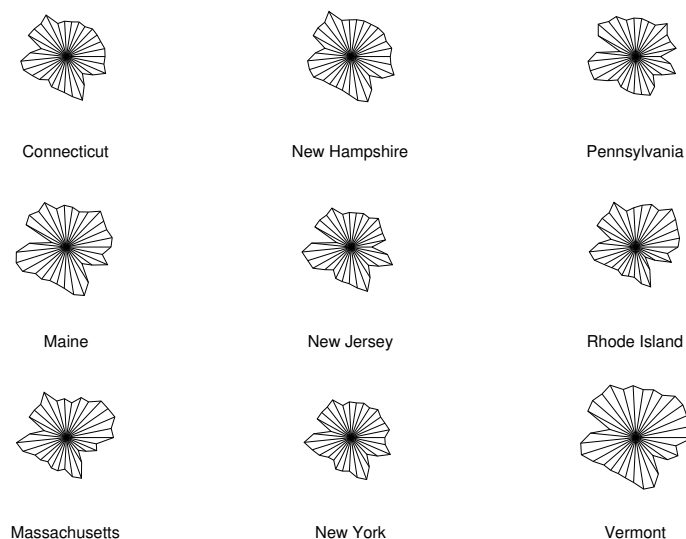


Abbildung 1.5: Star-Plot der Wahlergebnisse 1856-1976 für 9 Bundesstaaten der USA.

Eine andere bekannte, aber manchmal sehr irreführende multivariate Grafik sind die *Chernoff-Faces*. Augen, Ohren, Mund, Nase, Gesichtsform, ..., repräsentieren die verschiedenen Variablen. Die Objekte werden dann als „Gesicht“ dargestellt, wobei der Wert jeder Variable die Form oder Größe des Merkmals angibt.

1.2 Notation

Eine Matrix \mathbf{X} der Ordnung $(n \times p)$ wird geschrieben als

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} .$$

Generell wird eine Matrix mit fettgedruckten Großbuchstaben bezeichnet. Eine Matrix mit Spaltenordnung 1 wird Spaltenvektor genannt. Daher ist

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}$$

ein Spaltenvektor mit n Komponenten. Spaltenvektoren werden mit fettgedruckten Kleinbuchstaben bezeichnet. Ein Spaltenvektor kann auch transponiert dargestellt werden:

$$\mathbf{x}^\top = (x_1, \dots, x_n) .$$

Die Spalten einer Matrix \mathbf{X} werden geschrieben als

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p$$

und die Zeilen (geschrieben als Spaltenvektoren)

$$\mathbf{x}_{1.}, \mathbf{x}_{2.}, \dots, \mathbf{x}_{n.} .$$

Diese Notation ist eher unüblich, und oft wird das selbe Symbol für Zeilen oder Spalten einer Matrix verwendet. Dies kann aber zu Verwirrungen führen, und daher scheint diese Lösung sinnvoll zu sein. Somit ist

$$\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p) = \begin{pmatrix} \mathbf{x}_{1.}^\top \\ \mathbf{x}_{2.}^\top \\ \vdots \\ \mathbf{x}_{n.}^\top \end{pmatrix}$$

$$\text{mit } \mathbf{x}_j = \begin{pmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{pmatrix} \quad \text{und} \quad \mathbf{x}_{i.} = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix} .$$

1.3 Beschreibende Statistik

Wie im univariaten (eindimensionalen) Fall möchte man auch für multivariate Daten einfache statistische Kenngrößen angeben. Das Stichprobenmittel gibt Aufschluß über die Lage bzw. das Zentrum der betrachteten Größe. Sinnvollerweise wird man dieses Lagemaß für jede betrachtete Variable (d.h. für jede Spalte einer Datenmatrix) gesondert ermitteln.

Seien $x_{11}, x_{21}, \dots, x_{n1}$ die n Messungen der ersten Variable. Dann ist das *arithmetische Mittel* dieser Messungen für die erste Variable

$$\bar{x}_1 = \frac{1}{n} \sum_{i=1}^n x_{i1} .$$

Auf analoge Weise kann das arithmetische Mittel für die anderen Variablen bestimmt werden. Man erhält somit $\bar{x}_1, \dots, \bar{x}_p$ für p Variablen der $(n \times p)$ -Datenmatrix \mathbf{X} . Das *multivariate (arithmetische) Mittel* ist dann der p -dimensionale Vektor $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_p)^\top$.

Eine weitere wichtige Kenngröße der beschreibenden Statistik ist die *Varianz*. Für unsere Matrix \mathbf{X} ist die Stichprobenvarianz der ersten Variable

$$s_{11} = \frac{1}{n-1} \sum_{i=1}^n (x_{i1} - \bar{x}_1)^2 .$$

Im allgemeinen kann die Stichprobenvarianz für die p Variablen definiert werden durch

$$s_{jj} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2$$

für $j = 1, \dots, p$. Nachdem auch ein Variationsmaß zwischen den einzelnen Variablen angegeben werden kann (Kovarianz), sollte man diese Streuungsmaße in Matrixform anschreiben. Wir ordnen daher die Stichprobenvarianzen s_{jj} entlang der Hauptdiagonale einer (Varianz-)Kovarianz-Matrix \mathbf{S} an. Die *Kovarianz* s_{jk} , die den linearen Zusammenhang zweier Variablen j und k angibt, ist definiert als

$$s_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) .$$

Es ist offensichtlich, dass die Kovarianzmatrix symmetrisch ist, da $s_{jk} = s_{kj}$ gilt.

Ein weiteres wichtiges Maß ist der Korrelationskoeffizient. Er ist ein standardisiertes Maß für den linearen Zusammenhang zweier Variablen und ist definiert durch

$$r_{jk} = \frac{s_{jk}}{\sqrt{s_{jj}}\sqrt{s_{kk}}} = \frac{\sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k)}{\sqrt{\sum_{i=1}^n (x_{ij} - \bar{x}_j)^2} \sqrt{\sum_{i=1}^n (x_{ik} - \bar{x}_k)^2}}$$

für $j = 1, \dots, p$ und $k = 1, \dots, p$. Klarerweise gilt $r_{jj} = 1$. Die Elemente r_{jk} können in der Stichproben-Korrelationsmatrix \mathbf{R} zusammengefasst werden.

1.4 Eigenwerte und Eigenvektoren

Sei Σ eine $(p \times p)$ -Matrix (quadratisch). Dann ist

$$q(a) = |\Sigma - a\mathbf{I}| \quad (1.1)$$

ein Polynom der Ordnung p in a . Die p -ten Wurzeln von $q(a)$, nämlich a_1, \dots, a_p , nennt man Eigenwerte von Σ .

Es gilt dann für $i = 1, \dots, p$

$$|\Sigma - a_i\mathbf{I}| = 0 \quad , \quad (1.2)$$

sodass $\Sigma - a_i\mathbf{I}$ singulär ist. Es gibt daher einen Vektor $\gamma_i \neq \mathbf{0}$ mit

$$\Sigma\gamma_i = a_i\gamma_i \quad \text{für } i = 1, \dots, p. \quad (1.3)$$

γ_i wird bezeichnet als (Rechts-) Eigenvektor von Σ zum Eigenwert a_i . Ein reeller Eigenvektor γ_i heißt standardisiert, wenn gilt: $\gamma_i^\top \gamma_i = 1$.

Sind ξ und ζ Eigenvektoren zum Eigenwert a_i , dann sind $\xi + \zeta$ und $\alpha\xi$ ($\alpha \in \mathbb{R}$) auch Eigenvektoren zu a_i .

Satz 1.4.1 *Sei \mathbf{C} eine nichtsinguläre $(p \times p)$ -Matrix. Dann gilt:*

$$|\Sigma - a_i\mathbf{I}| = |\mathbf{C}| |\Sigma - a_i\mathbf{C}^{-1}\mathbf{C}| |\mathbf{C}^{-1}| = |\mathbf{C}\Sigma\mathbf{C}^{-1} - a_i\mathbf{I}| \quad . \quad (1.4)$$

D.h. Σ und $\mathbf{C}\Sigma\mathbf{C}^{-1}$ haben dieselben Eigenwerte. Ist γ ein Eigenvektor von Σ zum Eigenwert a_i , dann ist

$$a_i\mathbf{C}\gamma = \mathbf{C}\Sigma\mathbf{C}^{-1}(\mathbf{C}\gamma) \quad . \quad (1.5)$$

Daraus folgt, dass $\mathbf{C}\gamma$ ein Eigenvektor von $\mathbf{C}\Sigma\mathbf{C}^{-1}$ zum Eigenwert a_i ist.

Satz 1.4.2 *Alle Eigenwerte einer symmetrischen $(p \times p)$ -Matrix Σ sind reell.*

Beweis: Sei $\gamma = \xi + i\zeta$, $a = b + ic$, $\gamma \neq \mathbf{0}$.

Durch Einsetzen in $\Sigma\gamma = a\gamma$, kann Real- und Imaginärteil berechnet werden.

$$\Sigma\xi = b\xi - c\zeta \quad , \quad \Sigma\zeta = c\xi + b\zeta \quad . \quad (1.6)$$

Premultiplikation mit ζ^\top bzw. ξ^\top und Subtraktion liefert $c = 0$, d.h. a ist reell. γ kann auch reell gewählt werden, d.h. $\zeta = \mathbf{0}$. \square

Satz 1.4.3 (Spektralzerlegungssatz oder Jordanscher Zerlegungssatz)

Jede symmetrische $(p \times p)$ -Matrix Σ kann zerlegt werden in

$$\Sigma = \mathbf{\Gamma}\mathbf{A}\mathbf{\Gamma}^\top = \sum_{i=1}^p a_i\gamma_i\gamma_i^\top \quad , \quad (1.7)$$

wobei $\mathbf{A} = \text{Diag}(a_1, \dots, a_p)$ eine Diagonalmatrix mit den Eigenwerten von Σ ist und $\mathbf{\Gamma}$ eine orthogonale Matrix ist (d.h. $\mathbf{\Gamma}^\top = \mathbf{\Gamma}^{-1}$), deren Spalten standardisierte Eigenvektoren γ_i von Σ sind.

Beweis: Seien $\gamma_1, \dots, \gamma_p$ orthonormale Vektoren mit

$$\Sigma \gamma_i = a_i \gamma_i \quad (1.8)$$

für Zahlen a_i ($i = 1, \dots, p$).

Dann ist

$$\gamma_i^\top \Sigma \gamma_j = a_j \gamma_i^\top \gamma_j = \begin{cases} a_i & \text{für } i = j \\ 0 & \text{für } i \neq j \end{cases}, \quad (1.9)$$

oder in Matrixform

$$\Gamma^\top \Sigma \Gamma = \mathbf{A} \quad (1.10)$$

Daraus folgt:

$$\Gamma(\Gamma^\top \Sigma \Gamma) \Gamma^{-1} = \Sigma = \Gamma \mathbf{A} \Gamma^\top \quad (1.11)$$

Aus Satz 1.4.1 folgt, dass Σ und \mathbf{A} die selben Eigenwerte haben, nämlich genau die Elemente von \mathbf{A} .

Nun muss noch eine Orthonormalbasis aus Eigenvektoren gefunden werden. Sind $a_i \neq a_j$ verschiedene Eigenwerte mit Eigenvektor $\xi + \zeta$, dann gilt:

$$a_i \xi^\top \zeta = \xi^\top a_j \zeta = \xi^\top \Sigma \zeta = \zeta^\top \Sigma \xi = a_j \zeta^\top \xi \implies \zeta^\top \xi = 0 \quad (1.12)$$

D.h. bei einer symmetrischen Matrix sind die Eigenvektoren, die zu verschiedenen Eigenwerten gehören, orthogonal.

Der Eigenraum einer Matrix Σ mit Eigenwert a ist definiert durch $\ker(\Sigma - a\mathbf{I})$. Angenommen Σ hat k verschiedene Eigenwerte mit zugehörigen Eigenräumen H_1, \dots, H_k der Dimensionen r_1, \dots, r_k und sei $r = \sum_{j=1}^k r_j$.

Da verschiedene Eigenräume orthogonal sind, existiert eine orthonormale Menge $\{(\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_{r_1}), (\mathbf{e}_{r_1+1}, \mathbf{e}_{r_1+2}, \dots, \mathbf{e}_{r_1+r_2}), \dots, (\mathbf{e}_{r_1+r_2+\dots+r_{k-1}+1}, \dots, \mathbf{e}_{r_1+r_2+\dots+r_k})\}$, bzw. in kompakterer Schreibweise $\{\mathbf{e}_1, \dots, \mathbf{e}_r\}$, von Vektoren, sodass die Vektoren mit Index

$$\sum_{i=1}^{j-1} r_i + 1, \sum_{i=1}^{j-1} r_i + 2, \dots, \sum_{i=1}^{j-1} r_i + r_j - 1, \sum_{i=1}^j r_i \quad (1.13)$$

eine Basis von H_j bilden. r_j ist kleiner oder gleich der Vielfachheit des entsprechenden Eigenwertes. Durch Umordnung der Eigenwerte λ_i (falls notwendig) gilt

$$\Sigma \mathbf{e}_i = a_i \mathbf{e}_i \quad \text{für } i = 1, \dots, r, \quad (1.14)$$

wobei $r \leq p$ ist. Sind alle Eigenwerte verschieden, ist $r = p$ und $\gamma_{(i)} = \mathbf{e}_i$.

Es ist noch zu zeigen, dass der Fall $r < p$ zu einem Widerspruch führt.

Seien o.B.d.A. alle Eigenwerte von Σ strikt positiv. (Falls nicht, ersetzen wir Σ durch $\Sigma + \alpha \mathbf{I}$ für positives α , da Σ und $\Sigma + \alpha \mathbf{I}$ die selben Eigenvektoren haben.) Sei

$$\mathbf{B} = \Sigma - \sum_{i=1}^r a_i \mathbf{e}_i \mathbf{e}_i^\top \quad (1.15)$$

Dann ist

$$\text{tr } \mathbf{B} = \text{tr } \Sigma - \sum_{i=1}^r a_i (\mathbf{e}_i^\top \mathbf{e}_i) = \sum_{i=r+1}^p a_i > 0, \quad (1.16)$$

da $r < p$. \mathbf{B} hat mindestens einen Eigenwert $\theta \neq 0$. Sei $\mathbf{x} \neq \mathbf{0}$ der zugehörige Eigenvektor, so gilt für $1 \leq j \leq r$:

$$\begin{aligned} \theta \mathbf{e}_j^\top \mathbf{x} &= \mathbf{e}_j^\top \mathbf{B} \mathbf{x} = \mathbf{e}_j^\top \left(\mathbf{\Sigma} - \sum_{i=1}^r a_i \mathbf{e}_i \mathbf{e}_i^\top \right) \mathbf{x} \\ &= \left(a_j \mathbf{e}_j^\top - \sum_{i=1}^r a_i (\mathbf{e}_j^\top \mathbf{e}_i) \mathbf{e}_i^\top \right) \mathbf{x} = (a_j \mathbf{e}_j^\top - a_j \mathbf{e}_j^\top) \mathbf{x} = 0 \quad . \end{aligned}$$

D.h. \mathbf{x} ist orthogonal zu \mathbf{e}_j für $j = 1, \dots, r$. Daher ist

$$\theta \mathbf{x} = \mathbf{B} \mathbf{x} = \left(\mathbf{\Sigma} - \sum_{i=1}^r a_i \mathbf{e}_i \mathbf{e}_i^\top \right) \mathbf{x} = \mathbf{\Sigma} \mathbf{x} - \sum_{i=1}^r a_i \mathbf{e}_i (\mathbf{e}_i^\top \mathbf{x}) = \mathbf{\Sigma} \mathbf{x} \quad , \quad (1.17)$$

sodass \mathbf{x} ein Eigenvektor von $\mathbf{\Sigma}$ ist. D.h. $\theta = a_i$ für ein i und \mathbf{x} ist Linearkombination gewisser Vektoren \mathbf{e}_i , was ein Widerspruch zur Orthogonalität von \mathbf{x} und \mathbf{e}_i ist. \square

1.5 Erwartungswert und Kovarianzmatrix

Sei \mathbf{X} eine $(p \times q)$ -Matrix von Zufallsgrößen, die mit dem Symbol $\mathbf{X} = [(x_{ij})]$ bezeichnet wird. Der Erwartungswert dieser Matrix ist definiert als

$$E(\mathbf{X}) = [(E(x_{ij}))] \quad , \quad (1.18)$$

der $(p \times q)$ -Matrix der Erwartungswerte $E(x_{ij})$. Aus der Linearität des Erwartungswert-Operators folgt für Datenmatrizen \mathbf{A} , \mathbf{B} und \mathbf{C}

$$E(\mathbf{A} \mathbf{X} \mathbf{B} + \mathbf{C}) = \mathbf{A} E(\mathbf{X}) \mathbf{B} + \mathbf{C} \quad . \quad (1.19)$$

Auch die Kovarianz kann auf den mehrdimensionalen Fall übertragen werden. Seien \mathbf{x} und \mathbf{y} zwei Zufallsvektoren, die nicht notwendigerweise die selbe Dimension haben. Dann gilt

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = [(Cov(x_i, y_j))] = E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{y} - \boldsymbol{\mu}_y)^\top] = E(\mathbf{x} \mathbf{y}^\top) - \boldsymbol{\mu}_x \boldsymbol{\mu}_y^\top \quad , \quad (1.20)$$

wobei $\boldsymbol{\mu}_x$ bzw. $\boldsymbol{\mu}_y$ die Erwartungswerte von \mathbf{x} bzw. \mathbf{y} sind. Wenn \mathbf{x} und \mathbf{y} statistisch unabhängig sind, gilt $\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbf{O}$. Für $\mathbf{x} = \mathbf{y}$ wird statt $\text{Cov}(\mathbf{x}, \mathbf{x})$ nur $\text{Cov}(\mathbf{x})$ geschrieben und es gilt

$$\text{Cov}(\mathbf{x}) = [(Cov(x_i, x_j))] = E[(\mathbf{x} - \boldsymbol{\mu}_x)(\mathbf{x} - \boldsymbol{\mu}_x)^\top] \quad . \quad (1.21)$$

Die so erhaltene Matrix enthält in der Diagonale ($i = j$) die Varianzen und wird deshalb oft mit *Varianz-Kovarianz-Matrix* bezeichnet. Für Datenmatrizen \mathbf{A} und \mathbf{B} gilt, wie leicht zu zeigen ist,

$$\text{Cov}(\mathbf{A} \mathbf{x}, \mathbf{B} \mathbf{y}) = \mathbf{A} \text{Cov}(\mathbf{x}, \mathbf{y}) \mathbf{B}^\top \quad , \quad (1.22)$$

woraus weiters folgt

$$\text{Cov}(\mathbf{A} \mathbf{x}) = \mathbf{A} \text{Cov}(\mathbf{x}) \mathbf{A}^\top \quad . \quad (1.23)$$

1.6 Die multivariate Normalverteilung

Die Dichte der multivariaten Normalverteilung wird erhalten, indem von der Dichte der univariaten Normalverteilung

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-[(x-\mu)/\sigma]^2/2} \quad -\infty < x < \infty \quad (1.24)$$

die Distanz $(x - \mu)/\sigma$ durch eine multivariate Distanz

$$(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) := \text{MD}^2(\mathbf{x}) \quad (1.25)$$

ersetzt wird. Ebenso wird die Konstante $\frac{1}{\sqrt{2\pi\sigma^2}}$ ersetzt durch eine allgemeinere Konstante, die das Volumen, das durch die Dichtefunktion beschrieben wird, auf 1 standardisiert. $\text{MD}(\mathbf{x})$ wird mit *Mahalanobis-Distanz* (MD) bezeichnet.

Sei also $\mathbf{x} = (x_1, \dots, x_p)^\top$ ein p -dimensionaler Zufallsvektor. Der Erwartungswert von \mathbf{x} sei $\boldsymbol{\mu}$, die Kovarianzmatrix sei $\boldsymbol{\Sigma} = [(\sigma_{ij})]$, die positiv definit ist (Schreibweise $\boldsymbol{\Sigma} \geq \mathbf{O}$).

Definition 1.6.1 *Mit obigen Bezeichnungen ist \mathbf{x} p -dimensional normalverteilt, wenn für die Dichtefunktion von \mathbf{x} gilt*

$$f(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (1.26)$$

mit $-\infty < x_i < \infty$ für $i = 1, \dots, p$. Als Notation wird dafür $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ verwendet.

Eigenschaften: Für $\mathbf{x} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ gilt:

- (a) $\mathbf{x} - \boldsymbol{\mu} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma})$
- (b) Sei \mathbf{A} eine $(q \times p)$ -Matrix, dann gilt $\mathbf{Ax} \sim N_q(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$.
- (c) Wird die MD gleich einer Konstanten c gesetzt, also

$$(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) := c^2, \quad (1.27)$$

so wird dadurch ein Ellipsoid im p -dimensionalen Raum mit Zentrum $\boldsymbol{\mu}$ beschrieben (siehe auch Kapitel Hauptkomponentenanalyse). Die Dichte der p -variaten Normalverteilung ist also auf solchen Ellipsoiden konstant. Es gilt außerdem, dass für normalverteiltes \mathbf{x} die Größe

$$(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \sim \chi_p^2, \quad (1.28)$$

also einer χ^2 -Verteilung mit p Freiheitsgraden folgt (p ist die Dimension von \mathbf{x}).

Die Größe $\sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$ ist das multivariate Analogon zur univariaten Quadratsumme $\sum_{i=1}^n (y_i - \bar{y})^2$. Dies führt zur **Wishart-Verteilung**, die folgendermaßen definiert ist:

Definition 1.6.2 Eine symmetrische $(p \times p)$ -Matrix \mathbf{W} mit Zufallsgrößen folgt einer Wishart-Verteilung, wenn \mathbf{W} mit unabhängigen p -dimensionalen identisch nach $N_p(\mathbf{0}, \Sigma)$ -verteilten Zufallsvektoren $\mathbf{x}_1, \dots, \mathbf{x}_n$ dargestellt werden kann als

$$\mathbf{W} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top = \mathbf{X}^\top \mathbf{X} \quad . \quad (1.29)$$

Man schreibt dafür $\mathbf{W} \sim W_p(\Sigma, n)$ oder kurz $\mathbf{W} \sim W_p$ und spricht von einer Wishart-Verteilung mit n Freiheitsgraden.

Eigenschaften:

(a) Sei $\mathbf{W} \sim W_p(\Sigma, n)$ und \mathbf{A} eine $(p \times q)$ -Matrix, dann ist

$$\mathbf{A}^\top \mathbf{W} \mathbf{A} \sim W_q(\mathbf{A}^\top \Sigma \mathbf{A}, n) \quad .$$

(b) $\Sigma^{-1/2} \mathbf{W} \Sigma^{-1/2} \sim W_p(\mathbf{I}, n)$

(c) Für $\mathbf{W} \sim W_p(\Sigma, n)$ und einen festen p -dimensionalen Vektor \mathbf{a} mit $\mathbf{a}^\top \Sigma \mathbf{a} \neq 0$ gilt

$$\frac{\mathbf{a}^\top \mathbf{W} \mathbf{a}}{\mathbf{a}^\top \Sigma \mathbf{a}} \sim \chi_n^2 \quad .$$

(d) Für $p = 1$ sind die Elemente von $\mathbf{x} = (x_1, \dots, x_n)^\top$ identisch verteilt nach $N_1(0, \sigma^2)$ und $\mathbf{x}^\top \mathbf{x} \sim W_1(\sigma^2, n)$. D.h. die Verteilung $W_1(\sigma^2, n)$ ist die gleiche Verteilung wie $\sigma^2 \chi_n^2$.

(e) Für die Diagonalelemente von \mathbf{W} gilt $w_{ii} \sim \sigma_i^2 \chi_n^2$.

(f) Wenn $\mathbf{W}_1 \sim W_p(\Sigma, n_1)$ und $\mathbf{W}_2 \sim W_p(\Sigma, n_2)$ gilt und \mathbf{W}_1 und \mathbf{W}_2 unabhängig sind, dann ist $\mathbf{W}_1 + \mathbf{W}_2 \sim W_p(\Sigma, n_1 + n_2)$.

Die Dichtefunktion der multivariaten Normalverteilung spielt in der multivariaten Statistik eine besonders große Rolle, weil die meisten Methoden bzw. Tests auf der Annahme basieren, dass die Daten einer multivariaten Normalverteilung entstammen. Eine Approximation der multivariaten Normalverteilung erhält man durch folgenden Satz:

Satz 1.6.1 (Zentraler Grenzverteilungssatz)

Für n voneinander unabhängige Beobachtungen einer Grundgesamtheit mit Mittel $\boldsymbol{\mu}$ und Kovarianzmatrix Σ (endlich, nicht singulär), gilt für großes $n - p$ approximativ:

$$\sqrt{n}(\bar{\mathbf{x}} - \boldsymbol{\mu}) \sim N_p(\mathbf{0}, \Sigma) \quad (1.30)$$

und

$$n(\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \Sigma^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}) \sim \chi_p^2 \quad . \quad (1.31)$$

Wie kann man nun überprüfen, ob die Daten einer multivariaten Normalverteilung entstammen? Ähnlich wie im univariaten Fall können Tests oder Grafiken zu einer Antwort führen.

Gängige statistische **Tests für multivariate Normalverteilung** sind der χ^2 -Anpassungstest (Conover, 1980), der Kolmogoroff-Smirnoff-Test (Smirnov, 1948; Afifi und Azen, 1979) und der Shapiro-Wilks-Test (Shapiro und Wilk, 1965). Letzterer gibt i.a. die zuverlässigeren Resultate. Diese Tests werden hier nicht näher behandelt, es sei auf die zitierte Literatur verwiesen.

Wir wollen durch ein **grafisches Verfahren** die multivariate Normalverteilung beurteilen. Beim univariaten Q-Q-Plot (Hazen, 1914) werden die Quantile der empirischen Verteilung den Quantilen der Normalverteilung gegenübergestellt. Die Erweiterung zum multivariaten Fall nennt man **χ^2 -Plot** oder Gamma-Plot (Easton und McCulloch, 1990), der folgendermaßen konstruiert wird:

- Man berechne für jede Beobachtung \mathbf{x}_i die quadrierten Mahalanobis-Distanzen

$$MD_i^2 = (\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{S}^{-1} (\mathbf{x}_i - \bar{\mathbf{x}}) \quad (1.32)$$

($i = 1, \dots, n$) und sortiere diese Werte $MD_{(1)}^2 \leq \dots \leq MD_{(n)}^2$.

- Man berechne $\chi_{p,i/n}^2$, also die Quantile i/n der χ_p^2 -Verteilung, und zeichne die Punktepaare $(\chi_{p,i/n}^2, MD_{(i)}^2)$ in eine Grafik.

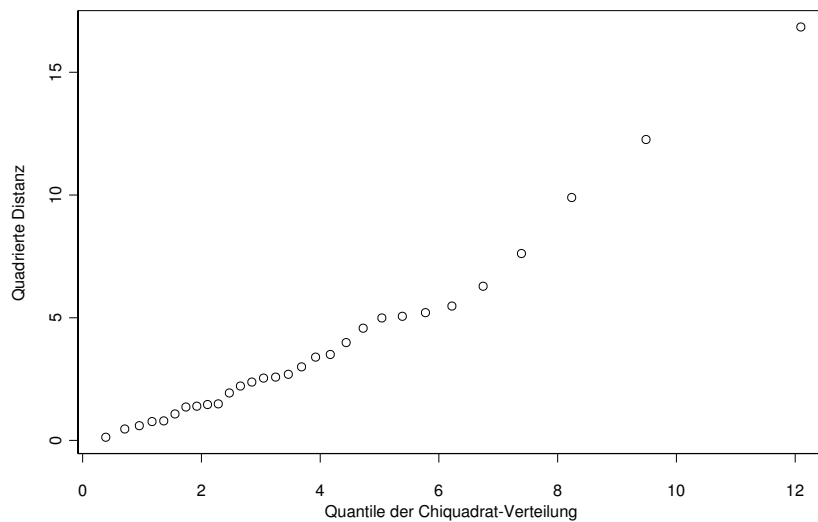
Liegen die Punkte in etwa auf einer Geraden, so kann man davon ausgehen, dass die Daten multivariat normalverteilt sind. Dieser grafische Test ist nur dann sinnvoll, wenn sowohl n als auch $n - p$ größer als 30 sind. Obwohl die Distanzen MD_i^2 nicht voneinander unabhängig sind oder exakt χ^2 -verteilt sind, ist diese Grafik sehr informativ.

Beispiel 1.6.1 *Wir betrachten die 4-dimensionalen Plattendaten aus Tabelle 1.2. Der entsprechende χ^2 -Plot ist in Abbildung 1.6 dargestellt. Man erkennt, dass die Ausreißer (Objekte 9 und 16) mit den größten Distanzen nicht gut ins Bild der multivariaten Normalverteilung passen.*

1.7 Transformation zu Normalverteilung

Wenn festgestellt wird, dass die vorliegenden Daten nicht normalverteilt sind, dann können durch Verfahren, die diese Voraussetzung benötigen, falsche Schlüsse gezogen werden. Es gibt aber noch die Möglichkeit, die Daten so zu transformieren, sodass sie einer Normalverteilung entsprechen. Man kann dann die üblichen Verfahren, die für normalverteilte Daten gültig sind, anwenden.

Eine Transformation der Daten ist nichts anderes als eine Skalenänderung. Sind z.B. die Daten extrem rechtsschief, kann eine Transformation mit dem Logarithmus Symmetrie um das Mittel sowie annähernde Normalverteilung bewirken.

Abbildung 1.6: χ^2 -Plot für die Plattendaten aus Tabelle 1.2.

Durch theoretische Überlegungen kann man je nach Datentyp folgende Transformationen empfehlen:

<i>Ursprüngliche Skalierung</i>	<i>Transformation</i>
Zähldaten y	\sqrt{y}
Verhältnisse, Anteile \hat{p}	$\text{logit}(\hat{p}) = \frac{1}{2} \log \left(\frac{\hat{p}}{1-\hat{p}} \right)$
Korrelationen r	Fisher's $z(r) = \frac{1}{2} \log \left(\frac{1+r}{1-r} \right)$

Oft kann man nur direkt anhand der Daten entscheiden, welche Transformation am besten geeignet ist. Eine Vielfalt von Möglichkeiten bietet die *Power-Transformation*: Sei x eine Beobachtung (nur positive Werte) und λ ein Parameter (der auch negativ sein kann). Dann ist die Power-Transformation definiert durch x^λ . Für $\lambda = 0$ wird $x^0 = \ln x$ definiert. Um die entsprechende Potenz (Power) zu finden, können das Histogramm oder ein Q-Q-Plot betrachtet werden. Ergibt sich gute Symmetrie um den Mittelwert und entspricht die Form der Verteilung am ehesten der Normalverteilung, so hat man die optimale Transformation gefunden.

Eine leicht modifizierte Transformation wird von Box und Cox (1964) vorgeschlagen:

$$x^{(\lambda)} = \begin{cases} \frac{x^\lambda - 1}{\lambda} & \lambda \neq 0 \\ \ln x & \lambda = 0 \end{cases} \quad (1.33)$$

Diese Transformation ist definiert für reelles λ und $x > 0$. Sind Beobachtungen x_1, \dots, x_n gegeben, ist die Box-Cox-Lösung für die richtige Wahl der Potenz λ die

Lösung, die den Ausdruck (Likelihood-Funktion)

$$l(\lambda) = -\frac{n}{2} \ln \left[\frac{1}{n} \sum_{i=1}^n (x_i^{(\lambda)} - \overline{x^{(\lambda)}})^2 \right] + (\lambda - 1) \sum_{i=1}^n \ln x_i \quad (1.34)$$

maximiert. Dabei ist $x_i^{(\lambda)}$ in (1.33) definiert, und

$$\overline{x^{(\lambda)}} = \frac{1}{n} \sum_{i=1}^n x_i^{(\lambda)} = \frac{1}{n} \sum_{i=1}^n \left(\frac{x_i^\lambda - 1}{\lambda} \right) \quad (1.35)$$

ist das arithmetische Mittel der transformierten Beobachtungen.

$l(\lambda)$ kann mittels Computer für viele Werte von λ berechnet werden, und mit einer Grafik der beiden Werte kann das Maximum leicht ermittelt werden. Meist entscheidet man sich dann für “schöne” Werte, d.h. Werte, die in der Nähe von $\hat{\lambda}$ liegen, wie z.B. $\lambda = 0$ (Logarithmus) oder $\lambda = \frac{1}{2}$ (Wurzel).

Obige Transformationen waren eigentlich für den univariaten Fall, also für *eine* Variable beschrieben. Im mehrdimensionalen Fall (für p Variablen) wird für jede einzelne Variable eine optimale Transformation durchgeführt. Entschließt man sich z.B. dazu, jede Variable der Box-Cox-Transformation zu unterziehen, so muss man die Potenzen $\lambda_1, \dots, \lambda_p$ durch maximieren einer Likelihoodfunktion für jedes einzelne λ_k suchen.

Es sei noch erwähnt, dass selbst die “beste” Transformation nicht notwendigerweise bewirken muss, dass die Daten nun tatsächlich normalverteilt sind. Eine solche Transformation bewirkt lediglich eine möglichst gute Annäherung an die Normalverteilung.

1.8 Tests, Konfidenzbereiche

Analog zum eindimensionalen Fall können auch bei multivariaten Daten statistische Tests durchgeführt werden. Bekannte Tests wie z.B. Tests auf den Mittelwert können leicht auf den multivariaten Fall übertragen werden. Der Unterschied ist nun, dass die p korrelierten Variablen *gemeinsam* analysiert werden.

Möchte man im univariaten Fall testen, ob das Mittel μ einer Grundgesamtheit gleich einem bestimmten Wert μ_0 ist, d.h. wir testen die Nullhypothese $H_0 : \mu = \mu_0$ gegen die Alternative $H_1 : \mu \neq \mu_0$, so lautet die entsprechende Teststatistik (bei unbekannter Varianz)

$$t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}}, \quad (1.36)$$

die t -verteilt ist mit $n - 1$ Freiheitsgraden (unter der Voraussetzung, dass eine Stichprobe einer normalverteilten Grundgesamtheit vorliegt).

Wir betrachten nun im multivariaten Fall eine Stichprobe einer normalverteilten Grundgesamtheit, d.h. $\mathbf{x}_1, \dots, \mathbf{x}_n \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Ein Test $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$ kann durchgeführt werden mit der Teststatistik

$$T^2 = (\bar{\mathbf{x}} - \boldsymbol{\mu}_0)^\top \left(\frac{1}{n} \mathbf{S} \right)^{-1} (\bar{\mathbf{x}} - \boldsymbol{\mu}_0), \quad (1.37)$$

wobei

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top \text{ und } \boldsymbol{\mu}_0 = (\mu_{10}, \dots, \mu_{p0})^\top.$$

Diese Teststatistik wird mit *Hotelling's T^2* bezeichnet (nach dem berühmten Statistiker Harold Hotelling), sie ist verteilt nach

$$T^2 \sim \frac{(n-1)p}{(n-p)} F_{p, n-p}, \quad (1.38)$$

also F-verteilt mit p und $n-p$ Freiheitsgraden. Somit gilt für ein Fehlniveau α :

$$\alpha = P \left[T^2 > \frac{(n-1)p}{(n-p)} F_{p, n-p; 1-\alpha} \right] \quad (1.39)$$

Im univariaten Fall kann zu einem Test zum Niveau α unmittelbar ein *Konfidenzintervall* konstruiert werden. Analoges gilt auch im multivariaten Fall, nur dass man hier einen *Konfidenzbereich* erhält. Anhand des obigen Tests kann man auf einfache Weise einen $100(1-\alpha)\%$ Konfidenzbereich für das Mittel einer p -dimensionalen normalverteilten Größe definieren. Es ist dies ein Ellipsoid, das bestimmt ist für alle $\boldsymbol{\mu}$, die

$$n(\bar{\mathbf{x}} - \boldsymbol{\mu})^\top \mathbf{S}^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}) \leq \frac{p(n-1)}{(n-p)} F_{p, n-p; 1-\alpha} \quad (1.40)$$

erfüllen.

Literatur

- A.A. Afifi and S.P. Azen. *Statistical Analysis. A Computer Oriented Approach*. Acad. Press, New York, 1979.
- T.W. Anderson. Asymptotic theory for principal component analysis. *Ann. Math. Statist.*, 34:122–148, 1963.
- R.A. Becker, W.S. Cleveland, and A.R. Wilks. Dynamic graphics for data analysis. *Statistical Science*, 2(4):355–395, 1987.
- G.E.P. Box and D.R. Cox. An analysis of transformations (with discussion). *J. Roy. Statist. Soc. B*, 26(2):211–252, 1964.
- W.R. Dillon and M. Goldstein. *Multivariate Analysis: Methods and Applications*. John Wiley & Sons, New York, 1984.
- G.S. Easton and R.E. McCulloch. A multivariate generalization of quantile-quantile plots. *Journal of the American Statistical Association*, 85(410):376–386, 1990.
- B. Everitt. *Graphical Techniques for Multivariate Data*. North-Holland, New York, 1978.
- D.M. Hawkins, editor. *Topics in Applied Multivariate Analysis*. Cambridge University Press, Cambridge, 1982.
- A. Hazen. Storage to be provided in impounding reservoirs for municipal water supply. *Transactions of the American Society of Civil Engineers*, 77:1529–1669, 1914.

- R. Johnson and D. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, London, 4th edition, 1998.
- A.M. Kshirsagar. *Multivariate Analysis*. M. Dekker, New York, 1972.
- K.V. Mardia, J.T. Kent, and J.M. Bibby. *Multivariate Analysis*. Acad. Press, London, 1979.
- P.J. Rousseeuw, I. Ruts, and J.W. Tukey. The bagplot: A bivariate boxplot. *The American Statistician*, 53(4):382–387, 1999.
- G.A.F. Seber. *Multivariate observations*. John Wiley & Sons, New York, 1984.
- S.S. Shapiro and M.B. Wilk. An analysis of variance test for normality. *Biometrika*, 52(3):591–611, 1965.
- N.V. Smirnov. Table for estimating the goodness of fit of empirical distributions. *Annals of Mathematical Statistics*, 19:279–281, 1948.
- J.W. Tukey. *Exploratory Data Analysis*. Addison-Wesley, Reading, Mass., 1977.

Kapitel 2

Clusteranalyse

2.1 Einleitung

Die Clusteranalyse bietet die Möglichkeit, Strukturen (Klassen) in einer Menge von Objekten (ev. auch Variablen) zu erkennen. Nimmt man an, dass eine Menge von n Objekten derart strukturiert ist, dass sie in mehrere Gruppen (Cluster, Klassen) zerfällt, so lassen sich mittels der Clusteranalyse diese Gruppen festlegen. Die Gruppenzugehörigkeiten der Objekte sollen dabei wesentlich durch den Grad der „Ähnlichkeit“ der Objekte festgelegt werden, d.h. Objekte eines Clusters sollen einander ähnlich sein (Homogenität innerhalb der Klassen), Objekte verschiedener Cluster sollen einander unähnlich sein (Heterogenität zwischen den Klassen).

Als *Distanz- oder Unähnlichkeitsmaß* zwischen dem i -ten und j -ten Objekt kann z.B. die *Euklidische Distanz*

$$d(i, j) = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2} = \|\mathbf{x}_i - \mathbf{x}_j\|_2 \quad (2.1)$$

oder die *Manhattan-Distanz*

$$d(i, j) = \sum_{k=1}^p |x_{ik} - x_{jk}| = \|\mathbf{x}_i - \mathbf{x}_j\|_1 \quad (2.2)$$

gewählt werden. Ausgehend von einer Distanzmatrix können verschiedene Clustermethoden angewandt werden.

2.2 Klassifikationstyp

Es kommen im wesentlichen vier Klassifikationstypen in Frage: Überdeckung, Partition, Quasihierarchie und Hierarchie. Bei allen vier Typen, die in Folge beschrieben werden, kann man zwischen exhaustiver und nicht-exhaustiver Klassifikation unterscheiden. Bei einer *exhaustiven Klassifikation* wird jedes Objekt klassifiziert, bei einer *nicht-exhaustiven Klassifikation* können auch einzelne Objekte nicht den Klassen zugeteilt werden.

Eine **Überdeckung** ist eine Klassifikation, bei der sich einzelne Klassen überschneiden dürfen. Es darf dabei aber nicht eine Klasse vollständig in einer anderen enthalten sein.

Eine **Partition** ist eine spezielle Überdeckung, bei der verlangt wird, dass kein Objekt zu mehr als einer Klasse gehört. D.h. die Klassen dürfen sich nicht überschneiden und sind somit disjunkt.

Nachdem Überdeckungen und Partitionen die Struktur einer Objektmenge oft nur sehr grob wiedergeben, wählt man manchmal die feineren Klassifikationstypen Quasihierarchie und Hierarchie.

Eine **Quasihierarchie** ist eine Klassifikation, die durch eine Folge von Überdeckungen gebildet wird. Man erhält somit eine Art „Stammbaum“ mit einander überschneidenden Klassen, dessen unterste Stufe die größte und dessen oberste Stufe die feinste Überdeckung der Quasihierarchie ist.

Eine **Hierarchie** ist analog zur Quasihierarchie definiert, nur müssen hier die Klassen disjunkt sein. Die Hierarchie wird somit durch eine Folge von Partitionen bestimmt. Mit Hilfe eines „Stammbaums“ können die einzelnen Stufen der Partitionierung dargestellt werden. Eine andere Möglichkeit zur grafischen Veranschaulichung bietet das *Dendrogramm* (siehe später).

2.3 Bewertungskriterien für Klassifikationen

Bei den verschiedenen Verfahren der Clusteranalyse möchte man hohe Homogenität (Gleichartigkeit) innerhalb der Klassen, gleichzeitig aber hohe Heterogenität (Verschiedenheit) zwischen den Klassen erreichen. Aus diesem Grund werden verschiedene Gütemaße für Homogenität und Heterogenität angegeben.

2.3.1 Maße für die Homogenität einer Klasse

Gegeben sei eine Klasse K_l ($l = 1, \dots, k$) aus einer Klassifikation κ . Die Homogenität dieser Klasse wird mit einem Homogenitätsmaß $h(K_l) \geq 0$ gemessen. Die resultierende Zahl wird umso kleiner, je homogener die Klasse K_l ist.

Die Homogenität einer Klasse kann mit Hilfe einer *Distanzmatrix* gemessen werden. Eine Distanzmatrix enthält die Distanzen oder Unähnlichkeiten (engl. *dissimilarities*) aller Paare von Objekten (siehe oben).

- Summe der Distanzen:

$$h(K_l) = \frac{1}{|K_l|} \sum_{\substack{i < j \\ i, j \in K_l}} d(i, j)$$

Die Normierungskonstante $|K_l|$ bezeichnet die Anzahl der Objekte in der Klasse K_l .

- Maximum der Distanzen:

$$h(K_l) = \max_{i, j \in K_l} d(i, j)$$

Dieses Maß entspricht den beiden unähnlichsten Objekten einer Klasse und ist somit relativ streng. Große Klassen werden schlecht beurteilt.

- Minimum der Distanzen:

$$h(K_l) = \min_{i,j \in K_l} d(i,j)$$

Große heterogene Klassen haben mit diesem Maß oft einen sehr kleinen Wert.

- Summe der Varianzen der Variablen:

Hier geht man nicht von der Distanzmatrix, sondern direkt von der $(n \times p)$ -Datenmatrix $\mathbf{X} = [(x_{ij})]$ aus.

Das Mittel der j -ten Variable ($j = 1, \dots, p$) in der Klasse K_l ist definiert als

$$\bar{x}_j = \frac{1}{|K_l|} \sum_{i \in K_l} x_{ij} .$$

Somit ist die Varianz der j -ten Variable in der Klasse K_l

$$s_j^2(K_l) = \frac{1}{|K_l| - 1} \sum_{i \in K_l} (x_{ij} - \bar{x}_j)^2 .$$

Als Homogenitätsmaß nimmt man die Summe dieser Varianzen

$$h(K_l) = \sum_{j=1}^p s_j^2(K_l) .$$

Für einelementige Mengen ist dieses Maß nicht verwendbar.

2.3.2 Maße für die Heterogenität zwischen den Klassen

Gegeben seien zwei Klassen K_{l_1} und K_{l_2} aus einer Klassifikation. Ein Heterogenitätsmaß $v(K_{l_1}, K_{l_2})$ soll nichtnegative Werte annehmen, die kleiner sind, je ähnlicher sich die beiden Klassen sind. Man verlangt außerdem Symmetrie, also $v(K_{l_1}, K_{l_2}) = v(K_{l_2}, K_{l_1})$.

Bei den Heterogenitätsmaßen muss man unterscheiden, ob die Klassen disjunkt sind oder ob sie sich überschneiden. Das Maß hängt also vom Klassifikationstyp ab.

Disjunkte Klassen:

Wir gehen wiederum von einer Distanzmatrix aus. Bekannte Maße für die Heterogenität zweier disjunkter Klassen K_{l_1} und K_{l_2} sind:

- Complete Linkage:

$$v(K_{l_1}, K_{l_2}) = \max_{i \in K_{l_1}, j \in K_{l_2}} d(i,j)$$

Die Verschiedenheit wird durch das unähnlichste Objektpaar bestimmt. Man erhält dadurch oft zu wenig Cluster.

- Single Linkage:

$$v(K_{l_1}, K_{l_2}) = \min_{i \in K_{l_1}, j \in K_{l_2}} d(i, j)$$

Die Verschiedenheit wird durch das ähnlichste Objektpaar bestimmt. Man erhält dadurch oft zu viele Cluster.

- Average Linkage:

$$v(K_{l_1}, K_{l_2}) = \frac{1}{|K_{l_1}| |K_{l_2}|} \sum_{i \in K_{l_1}} \sum_{j \in K_{l_2}} d(i, j)$$

Die Verschiedenheit wird durch die durchschnittliche Ähnlichkeit der beiden Klassen bestimmt. Dieses Verfahren ist somit ein Kompromiss zwischen den beiden oben beschriebenen Methoden.

- Centroid Methode:

Hier geht man nicht von der Distanzmatrix, sondern direkt von der $(n \times p)$ -Datenmatrix $\mathbf{X} = [(x_{ij})]$ aus.

Man berechnet den euklidischen Abstand zwischen den Zentren (centroid), die definiert sind als

$$\bar{\mathbf{x}}(K_{l_m}) = \frac{1}{|K_{l_m}|} \left(\sum_{i \in K_{l_m}} x_{i1}, \dots, \sum_{i \in K_{l_m}} x_{ip} \right)^{\top} \quad \text{für } m = 1, 2,$$

das entspricht den Mittelwertvektoren für die p Variablen in den Klassen K_{l_1} und K_{l_2} . Das Heterogenitätsmaß ist dann

$$v(K_{l_1}, K_{l_2}) = \|\bar{\mathbf{x}}(K_{l_1}) - \bar{\mathbf{x}}(K_{l_2})\|_2.$$

Überschneidende Klassen:

Obige Heterogenitätsmaße werden so modifiziert, dass die Verschiedenheit der um die identischen Objekte reduzierten Klassen berechnet wird.

Ist eine Klasse vollständig in einer anderen enthalten (Hierarchie, Quasihierarchie), so wird die größere Klasse um die in ihr enthaltenen reduziert.

Beispiel: $K_{l_1} = \{8, 10\}$, $K_{l_2} = \{2, 7, 8, 10, 12\} \implies K_{l_2}^* = K_{l_2} - K_{l_1} \cap K_{l_2} = \{2, 7, 12\} \implies v(K_{l_1}, K_{l_2}^*)$

2.3.3 Maße für die Güte einer Klassifikation

Durch Maße für Homogenität und Heterogenität werden die einzelnen Klassen bzw. die Klassenpaare einer Klassifikation bewertet. Mit einem Gütemaß möchte man die Klassifikation selbst bewerten.

Das Gütemaß wird umso kleiner, je „besser“ die Klassifikation ist. Dies hängt aber wiederum vom Klassifikationstyp ab.

Als Gütemaß für **Überdeckungen** könnte man die Summe aller Klassenhomogenitäten minus der Summe aller Klassenhomogenitäten aller Teilmengen von Objekten, die gleichzeitig mehreren Klassen angehören, heranziehen. Nachdem dieses Maß nur von den Homogenitäten abhängt, könnte man diese Größe noch durch ein Maß für die Heterogenitäten dividieren.

Bei **Partitionen** könnte man die Summe der Klassenhomogenitäten als Gütemaß verwenden. Eine andere Möglichkeit ist ein normierter Kehrwert der Heterogenitäten oder der Quotient von Homogenitäten und Heterogenitäten.

Gütemaße für **Hierarchien** und **Quasihierarchien** werden für jede Stufe in der (Quasi-)Hierarchie berechnet, wobei obige Maße verwendet werden können.

2.4 Konstruktionsverfahren

Die Verfahren zur Konstruktion einer Clusterung hängen natürlich vom Klassifikationstyp ab. Aus Platzgründen wollen wir in Folge nur Verfahren für die gängigeren Typen Partition und Hierarchie betrachten.

2.4.1 Partitionen (Partitionierungsmethoden)

Ein iteratives Verfahren für große Objektmenen

Die Klassenzahl k muss vorgegeben werden, wobei man dieses Verfahren ev. mit verschiedener Klassenzahl wiederholen kann. Man wählt eine Anfangspartition, d.h. k Objekte werden zufällig ausgewählt (Zentralobjekte). Die restlichen $n - k$ Objekte werden jeweils dem Zentralobjekt zugeordnet, dem sie am ähnlichsten sind (kleinste Distanz). Die Anfangspartition ist jetzt komplett.

Dann wird für jedes Objekt, das nicht zu einer einelementigen Klasse gehört, geprüft, ob die Güte verkleinert wird, wenn das Objekt in eine der anderen Klassen transferiert wird. Ist dies möglich, so wird jenes Objekt transferiert, das die größte Abnahme der Güte bewirkt. Es wird so lange iteriert, bis keine wesentliche Verbesserung der Güte möglich ist. Diese Methode kann ev. mit verschiedenen Anfangspartitionen durchgeführt werden.

Bei der **k-means-Methode** werden, ausgehend von der Anfangspartition, so lange die Objekte in andere Cluster transferiert, bis ein „Fehlermaß“ nicht mehr reduziert werden kann. Ein solches Fehlermaß ist z.B. die Summe über alle Objekte von der quadrierten euklidischen Distanz jedes Punktes zu seinem Clustercentroid, das jedesmal neu berechnet werden muss.

Ein rekursives Verfahren für große Objektmenen

Hier werden die Klassen nacheinander bestimmt, die Klassenanzahl braucht somit nicht vorgegeben werden.

Als Zentralobjekt der ersten Klasse wird eines jener beiden Objekte mit minimalem Abstand zueinander gewählt. Alternativ kann jenes Objekt i mit größter „Punktdichte“ gewählt werden, d.h. die Anzahl der Objekte j mit $d(i, j) \leq \tilde{d}$ soll maximiert werden. Ausgehend von diesem Kernobjekt wird jenes Objekt in die Klasse

genommen, das minimalen Abstand (größte Ähnlichkeit) zum Kernobjekt besitzt. Dies wird solange fortgeführt, bis eine Homogenitätsschranke überschritten wird oder die Homogenität von einem Schritt zum nächsten radikal vergrößert wird. Aus der Restmenge wird, analog zu oben, wieder ein Kernobjekt gewählt und Objekte adjungiert. Das Verfahren wird abgebrochen, wenn alle Objekte klassifiziert sind (exhaustiv) oder wenn nur noch sehr kleine Gruppen auftreten.

2.4.2 Hierarchie

Hierarchische Verfahren werden durch eine Folge von Partitionen (disjunkte Klassen) gebildet. Die Konstruktion kann divisiv (von einer groben Partition schrittweise immer feinere erzeugen) oder agglomerativ erfolgen. Divisive Verfahren sind sehr rechenaufwendig und daher wenig gebräuchlich.

Agglomerative Methoden:

Ausgehend von n einelementigen Klassen (Anfangspartition) werden in jedem Schritt jene beiden Klassen zusammengefasst, die minimal verschieden sind, bis nur noch eine Klasse übrigbleibt. Zur Feststellung, welche Klassen größte Ähnlichkeit (kleinste Verschiedenheit) haben, wählt man ein Distanzmaß (siehe Maße für Heterogenität disjunkter Klassen) wie z.B. **single linkage**, **complete linkage**, **average linkage**, **centroid method**. Jene Klassen mit minimalem Distanzmaß werden vereint.

Single linkage erzeugt viele (kleine) Klassen (Ketteneffekt) und ermöglicht eine Erkennung stark verzweigter Klassen. Nachteil ist, dass oft zwei unterschiedliche Klassen nur dadurch vereint werden, weil zwei Objekte dieser Klassen geringen Abstand haben. Große Klassen werden früh vereint.

Complete linkage führt bei inhomogenen Klassen zu schlechten Gruppierungen. Man erhält oft wenige relativ große Klassen.

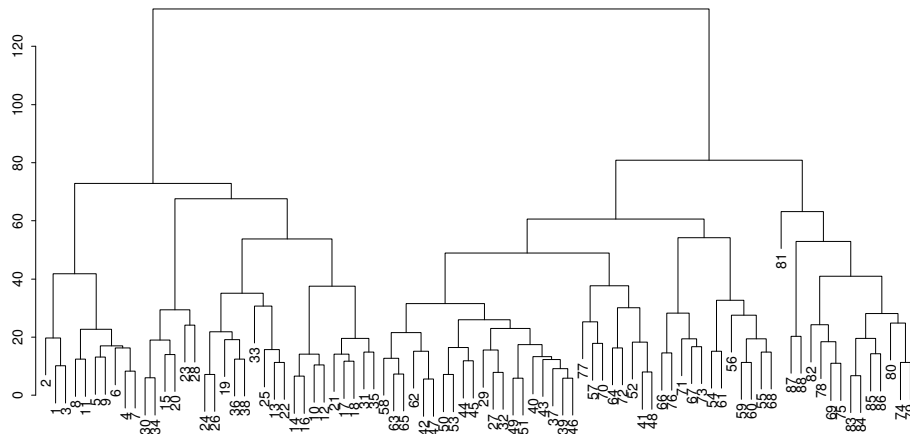
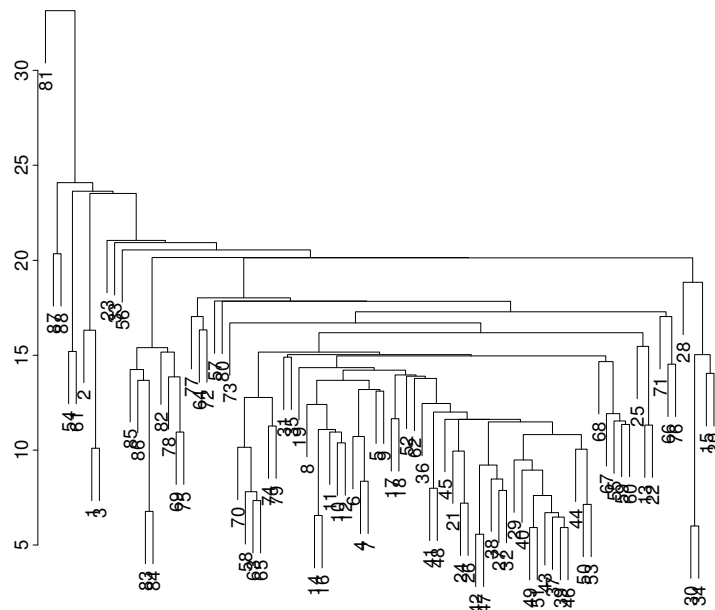
Agglomerative Verfahren können beendet werden, wenn nur noch eine Gruppe übrig bleibt, oder wenn eine Heterogenitätsschranke überschritten wird.

Da die Distanzen der Gruppen in jedem Schritt größer werden (Monotonie), können die Partitionen der Hierarchie grafisch in einem **Dendrogramm** dargestellt werden. Bei Fusionierung zweier Klassen kann das Distanzmaß (Heterogenitätsmaß) angegeben werden, das als Gütemaß dient.

Beispiel 2.4.1 Von 88 Studenten wurden die Prüfungsergebnisse in den Fächern *Mechanik*, *Analytische Geometrie*, *Lineare Algebra*, *Analysis* und *Elementare Statistik* aufgezeichnet. Von den 100 erreichbaren Punkten gab es die in Tabelle 2.1 angeführten Ergebnisse. Die 88 Studenten sollen anhand ihrer Leistungen in den fünf Prüfungsfächern geclustert werden. Die Clusterung wird mit *complete linkage* sowie mit *single linkage* durchgeführt. Die Resultate sind mit Hilfe von Dendrogrammen in Abbildung 2.1 und 2.2 dargestellt.

Tabelle 2.1: Prüfungsergebnisse von 88 Studenten in den Fächern Mechanik (ME), Analytische Geometrie (AG), Lineare Algebra (LA), Analysis (AN) und Elementare Statistik (ES); 100 Punkte waren erreichbar.

Student	ME	AG	LA	AN	ES	Student	ME	AG	LA	AN	ES
1	77	82	67	67	81	45	46	61	46	38	41
2	63	78	80	70	81	46	40	57	51	52	31
3	75	73	71	66	81	47	49	49	45	48	39
4	55	72	63	70	68	48	22	58	53	56	41
5	63	63	65	70	63	49	35	60	47	54	33
6	53	61	72	64	73	50	48	56	49	42	32
7	51	67	65	65	68	51	31	57	50	54	34
8	59	70	68	62	56	52	17	53	57	43	51
9	62	60	58	62	70	53	49	57	47	39	26
10	64	72	60	62	45	54	59	50	47	15	46
11	52	64	60	63	54	55	37	56	49	28	45
12	55	67	59	62	44	56	40	43	48	21	61
13	50	50	64	55	63	57	35	35	41	51	50
14	65	63	58	56	37	58	38	44	54	47	24
15	31	55	60	57	73	59	43	43	38	34	49
16	60	64	56	54	40	60	39	46	46	32	43
17	44	69	53	53	53	61	62	44	36	22	42
18	42	69	61	55	45	62	48	38	41	44	33
19	62	46	61	57	45	63	34	42	50	47	29
20	31	49	62	63	62	64	18	51	40	56	30
21	44	61	52	62	46	65	35	36	46	48	29
22	49	41	61	49	64	66	59	53	37	22	19
23	12	58	61	63	67	67	41	41	43	30	33
24	49	53	49	62	47	68	31	52	37	27	40
25	54	49	56	47	53	69	17	51	52	35	31
26	54	53	46	59	44	70	34	30	50	47	36
27	44	56	55	61	36	71	46	40	47	29	17
28	18	44	50	57	81	72	10	46	36	47	39
29	46	52	65	50	35	73	46	37	45	15	30
30	32	45	49	57	64	74	30	34	43	46	18
31	30	69	50	52	45	75	13	51	50	25	31
32	46	49	53	59	37	76	49	50	38	23	9
33	40	27	54	61	61	77	18	32	31	45	40
34	31	42	48	54	68	78	8	42	48	26	40
35	36	59	51	45	51	79	23	38	36	48	15
36	56	40	56	54	35	80	30	24	43	33	25
37	46	56	57	49	32	81	3	9	51	47	40
38	45	42	55	56	40	82	7	51	43	17	22
39	42	60	54	49	33	83	15	40	43	23	18
40	40	63	53	54	25	84	15	38	39	28	17
41	23	55	59	53	44	85	5	30	44	36	18
42	48	48	49	51	37	86	12	30	32	35	21
43	41	63	49	46	34	87	5	26	15	20	20
44	46	52	53	41	40	88	0	40	21	9	14

Abbildung 2.1: Clusterung der Prüfungsdaten (*Complete Linkage*)Abbildung 2.2: Clusterung der Prüfungsdaten (*Single Linkage*)

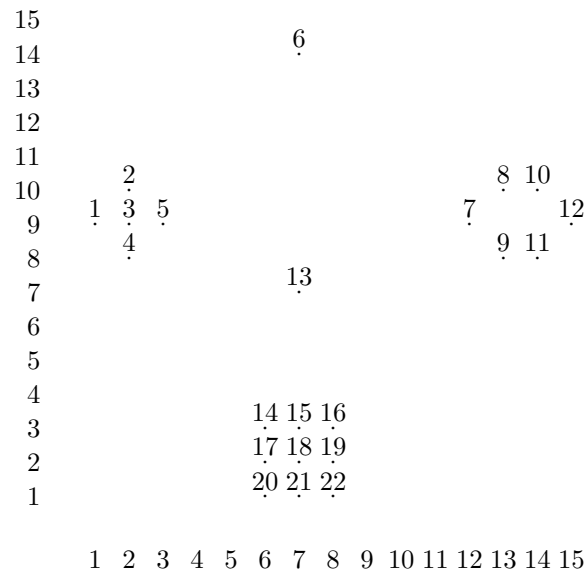
2.5 Fuzzy Clusterung

Bei den bisher besprochenen Methoden wurde jedes Objekt (exhaustiv) genau einem Cluster zugeordnet (man spricht von *harder* Clusterung). Bei der *Fuzzy*-Clusterung wird jedes Objekt auf alle Cluster „prozentuell“ aufgeteilt, wobei diese Aufteilung je nach Ähnlichkeit zum Cluster erfolgt. Die Anzahl der Cluster muss hier vorgegeben

werden.

Daten wie in Abbildung 2.3 können durch die Fuzzy-Clusterung wesentlich besser beschrieben werden als durch Partitionierungsmethoden.

Abbildung 2.3: Datensatz mit zwei Objekten, die nicht eindeutig einem bestimmten Cluster zugeordnet werden können



Dieses Beispiel lässt deutlich drei Cluster erkennen sowie die Objekte 6 und 13, die jedoch nicht eindeutig einem Cluster zugeordnet werden können.

Diese Methode kann aussagen, dass z.B. Objekt 1 fast sicher zu Cluster 1 gehört, wohingegen Objekt 13 mit praktisch gleicher Chance einem der drei Cluster zugeteilt werden sollte.

Jener Koeffizient, der angibt, zu welchem Teil ein Objekt zu einem bestimmten Cluster gehört, hat einen Wert zwischen 0 und 1 (0-100%) und heißt *Zugehörigkeitskoeffizient*.

Die Zugehörigkeitskoeffizienten zu den Daten aus Abbildung 2.3 sind in Tabelle 2.2 dargestellt. Wie aus diesem Beispiel auch ersichtlich ist, beträgt die Summe der Zugehörigkeitskoeffizienten eines Objekts 1 (100%). Der wesentliche Unterschied dieser Methode zur harten Clusterung liegt darin, dass sie detailliertere Information über die Struktur der Daten liefert. Der Nachteil dabei ist, dass mit wachsender Anzahl der Objekte der Ausdruck der Ergebnisse unüberschaubar wird und so die erhaltene Information oft nicht verarbeitet werden kann.

Ein möglicher Algorithmus wäre folgender: Es gibt keine repräsentativen Objekte, stattdessen wird versucht, die Objektfunktion

$$\sum_{v=1}^k \frac{\sum_{i,j=1}^n u_{iv}^2 u_{jv}^2 d(i,j)}{2 \sum_{j=1}^n u_{jv}^2} \quad (2.3)$$

Tabelle 2.2: Zugehörigkeitskoeffizienten zu den Daten in Abbildung 6.2

Objekt	Cluster 1	Cluster 2	Cluster 3
1	0.87	0.06	0.07
2	0.88	0.05	0.07
3	0.93	0.03	0.04
4	0.86	0.06	0.08
5	0.87	0.06	0.07
6	0.42	0.35	0.23
7	0.08	0.82	0.10
8	0.06	0.87	0.07
9	0.06	0.86	0.08
10	0.06	0.87	0.07
11	0.06	0.86	0.08
12	0.07	0.84	0.09
13	0.36	0.27	0.37
14	0.12	0.08	0.80
15	0.08	0.07	0.85
16	0.10	0.10	0.80
17	0.08	0.06	0.86
18	0.04	0.04	0.92
19	0.07	0.07	0.86
20	0.10	0.08	0.82
21	0.07	0.06	0.87
22	0.09	0.09	0.82

zu minimieren. Sie enthält nur das Ähnlichkeitsmaß $d(i, j)$ und den gesuchten Zugehörigkeitskoeffizienten u_{iv} des Objekts i zum Cluster v . Die Summe im Zähler geht über alle Objektpaare $\{i, j\}$ (anstelle der Distanzsumme der Objekte zu den Clustermittelpunkten bei anderen Verfahren). Da doppelt aufsummiert wird ($\{i, j\} = \{j, i\}$), wird durch 2 dividiert. Die äußere Summe geht über alle Cluster v , wodurch die zu minimierende Objektfunktion ein Maß für die Streuung ist.

Verwandte Methoden

Als eine der ersten Fuzzy-Clustermethoden wurde die *fuzzy k-means*-Methode entwickelt (siehe Bezdek, 1974, Dunn, 1974). Diese Methode minimiert die Objektfunktion

$$\sum_i \sum_v u_{iv}^2 \|\mathbf{x}_i - \mathbf{m}_v\|_2^2 = \sum_i \sum_v u_{iv}^2 \sum_{j=1}^p (x_{ij} - m_{vj})^2 \quad . \quad (2.4)$$

Dabei ist \mathbf{m}_v das Zentrum des Clusters v , und \mathbf{x}_i bezeichnet das i -te Objekt. Für jede Variable (Index j) wird das Zentrum folgendermaßen berechnet:

$$m_{vj} = \frac{\sum_i u_{iv}^2 x_{ij}}{\sum_i u_{iv}^2} . \quad (2.5)$$

Die Norm $\|\mathbf{x}_i - \mathbf{m}_v\|_2$ ist die euklidische Distanz zwischen Objekt \mathbf{x}_i und dem Zentrum des Clusters v . Wesentlich bei dieser Methode ist die Annahme, dass die verschiedenen Objekte durch ihre Koordinaten im p -dimensionalen Raum gegeben sind. Das ist eine Einschränkung im Vergleich zu der oben beschriebenen Methode, die eine solche Forderung nicht benötigt, sondern nur Distanzen oder Unähnlichkeiten zwischen den Objekten voraussetzt. Die beiden Methoden können, falls die Daten aus Messungen bestehen, direkt verglichen werden. Durch Umformungen lässt sich zeigen, dass folgende Gleichheit gilt:

$$\sum_i \sum_v u_{iv}^2 \sum_j (x_{ij} - m_{vj})^2 = \sum_v \frac{\sum_i \sum_j u_{iv}^2 u_{jv}^2 \|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2 \sum_j u_{jv}^2} . \quad (2.6)$$

Der letzte Ausdruck ist genau die Objektfunktion der oben beschriebenen Methode, mit dem Unterschied, dass die Distanzen quadriert sind. Das veranschaulicht die Äquivalenzen zwischen der oben beschriebenen und der *fuzzy k-means*-Methode, wenn die Objekte durch Messungen beschrieben werden. (Als Eingabe für die erste Methode wird eine Distanzmatrix, in der die Quadrate der Euklidischen Abstände zwischen den Objekten eingetragen werden, benötigt.)

Eine andere Methode ist der *MND2-Algorithmus* (siehe Roubens, 1978). Dabei wird folgende Objektfunktion minimiert:

$$\sum_v \sum_{i,j} u_{iv}^2 u_{jv}^2 d(i,j) . \quad (2.7)$$

Auch dieser Ausdruck ist der Objektfunktion der ersten Methode ähnlich. Im Nenner der Objektfunktion der ersten Methode steht der Ausdruck $\sum_j u_{jv}^2$. Dieser Term ist der wesentliche Unterschied zwischen den beiden Objektfunktionen.

Literatur

- J.C. Bezdek. Cluster validity with fuzzy sets. *Cybernetics*, 3:58–72, 1974. Scripta Publ. Comp., Washington, D.C.
- H.H. Bock. Automatische Klassifikation. In K.P. Grotemeyer, D. Morgenstern, and H. Tietz, editors, *Studia Mathematica/Mathematische Lehrbücher*, volume XXIV. Vandenhoeck & Ruprecht, Göttingen, 1974.
- J.C. Dunn. A fuzzy relative of the isodata progress and its use in detecting compact well-separated clusters. *Cybernetics*, 3:32–57, 1974. Scripta Publ. Comp., Washington, D.C.
- J.A. Hartigan. *Clustering algorithms*. Wiley & Sons, New York, 1975.

- L. Kaufman and P.J. Rousseeuw. *Finding Groups in Data*. Wiley & Sons, New York, 1990.
- M. Roubens. Pattern classification problems and fuzzy sets. *Fuzzy Sets and Systems*, 1:239–253, 1978.

Kapitel 3

Multivariate lineare Regression

3.1 Einführung

Die multivariate Regressionsanalyse ist ein Instrumentarium zur Vorhersage von Werten einer oder mehrerer quantitativer Größen, den abhängigen Variablen, durch bekannte (unabhängige) Größen. Es wird versucht, einen funktionellen Zusammenhang (der hier als linear vorausgesetzt ist) zwischen den abhängigen und den unabhängigen Größen zu finden.

Wir werden zunächst den Fall betrachten, wo nur eine einzige abhängige Variable vorliegt. Später wird dieses Modell erweitert für mehrere abhängige Variablen.

3.2 Lineare multiple Regression

Seien x_1, \dots, x_q die q Variablen, die zur Vorhersage einer abhängigen Variable Y vorliegen. Z.B. könnte die abhängige Variable der Preis eines neuen Autos sein, die unabhängigen Größen könnten Hubraum, Leistung, Gewicht, Luxusategorie etc. sein. Das lineare multiple Regressionsmodell sagt aus, dass Y dargestellt werden kann durch ein Mittel, das kontinuierlich von x_j abhängt, sowie von einem zufälligen Fehler ε . Dieser Fehlerterm beinhaltet sowohl Messfehler als auch den Einfluss anderer Variablen, die in diesem Modell nicht berücksichtigt wurden. Die Werte der Vorhersagevariablen x_1, \dots, x_q werden als *fest* angesehen, der Fehlerterm und damit auch die abhängige Größe sind *zufällig*.

Das lineare Regressionsmodell mit *einer abhängigen Größe* hat somit die Form

$$Y = \beta_0 + \beta_1 x_1 + \dots + \beta_q x_q + \varepsilon . \quad (3.1)$$

Die unbekannten *Regressionsparameter* $\beta_0, \beta_1, \dots, \beta_q$ geben den funktionellen Zusammenhang an, der bei diesem Modell *linear* ist. Diese Regressionskoeffizienten können aufgrund von konkreten Beobachtungen in weiterer Folge geschätzt werden.

Sind nun n voneinander unabhängige Beobachtungen sowohl von Y als auch von den zugehörigen unabhängigen Größen x_1, \dots, x_q gegeben, so ist das gesamte Modell (**multiple Regressionsmodell**)

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_q x_{iq} + \varepsilon_i$$

$$\begin{aligned}
Y_2 &= \beta_0 + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_q x_{2q} + \varepsilon_2 \\
&\vdots \\
Y_n &= \beta_0 + \beta_1 x_{n1} + \beta_2 x_{n2} + \dots + \beta_q x_{nq} + \varepsilon_n
\end{aligned}$$

oder in Matrixschreibweise:

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} . \quad (3.2)$$

\mathbf{Y} ist ein n -dimensionaler Zufallsvektor. Um den Unterschied zu einem Vektor mit konkreten Beobachtungen (später mit \mathbf{y} bezeichnet) besser hervorzuheben, ist der Zufallsvektor \mathbf{Y} ausnahmsweise (und nur in diesem Kapitel) mit großem fettgedrucktem Symbol bezeichnet.

Die erste Spalte der *Design-Matrix* \mathbf{X} besteht aus Einsen, da sie Multiplikator der Konstanten β_0 ist. Somit ist \mathbf{X} eine $(n \times (q + 1))$ -Matrix.

Von den Fehlertermen wird folgendes vorausgesetzt:

1. $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ und
2. $\text{Cov}(\boldsymbol{\varepsilon}) = E(\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top) = \sigma^2 \mathbf{I}_n$.

σ^2 beschreibt hier die Varianz des Fehlerterms, die für alle Komponenten gleich ist. Außerdem sind verschiedene Komponenten des Fehlerterms unkorreliert.

Zum Testen von Hypothesen bzw. für die Konstruktion von Konfidenzintervallen wird man außerdem noch gemeinsame Normalverteilung voraussetzen müssen.

3.3 Der Kleinste-Quadratsummen-Schätzer

In der Regressionsanalyse wird nun versucht, aufgrund gegebener Werte der Vorhersage-Variablen den genauen funktionellen Zusammenhang zur abhängigen Variable festzustellen. Man versucht also, aufgrund von bekannten Werten $1, x_{i1}, \dots, x_{iq}$ das Regressionsmodell (3.2) den beobachteten Werten y_i anzupassen. Genauer gesagt, es sollen die Regressionskoeffizienten $\boldsymbol{\beta}$ und die Fehlervarianz σ^2 aufgrund der vorliegenden Daten ermittelt werden.

Die Differenz $y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_q x_{iq}$, die die Abweichung der Anpassung vom beobachteten Wert angibt, wird auch als *Residuum* bezeichnet. Bei der *Methode der Kleinsten Quadrate* soll die Summe der quadrierten Residuen minimiert werden:

$$\begin{aligned}
S(\boldsymbol{\beta}) &= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_q x_{iq})^2 \\
&= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})
\end{aligned} \quad (3.3)$$

Jener Koeffizient, der dieses Kriterium minimiert, wird mit $\hat{\boldsymbol{\beta}}$ bezeichnet. Er ist der Kleinste-Quadratsummen-Schätzer des Regressionsparameters $\boldsymbol{\beta}$.

Gleichung (3.3) kann ausmultipliziert werden:

$$S(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta} . \quad (3.4)$$

Da es sich bei den Summanden um Skalare handelt, gilt

$$\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} = (\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y})^\top = \mathbf{y}^\top \mathbf{X} \boldsymbol{\beta} \quad . \quad (3.5)$$

Die Residuen-Quadratsumme (QS) ist somit

$$S(\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{y} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} \quad . \quad (3.6)$$

Partielles Differenzieren nach dem Parametervektor $\boldsymbol{\beta}$ ergibt

$$\frac{\partial S(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0 - 2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \boldsymbol{\beta} \quad . \quad (3.7)$$

Das Minimum der Residuen-QS erhält man durch Nullsetzen obiger Gleichung. Der Kleinste-QS-Schätzer für $\boldsymbol{\beta}$ ist dann

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad . \quad (3.8)$$

Bemerkung: Hat \mathbf{X} nicht vollen Rang $q+1 \leq n$, so muss $(\mathbf{X}^\top \mathbf{X})^{-1}$ ersetzt werden durch eine Verallgemeinerte Inverse (z.B. Moore-Penrose) $(\mathbf{X}^\top \mathbf{X})^-$.

Mit Hilfe von $\hat{\boldsymbol{\beta}}$ können die *geschätzten Werte* $\hat{\mathbf{y}}$ von \mathbf{y} berechnet werden durch

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\beta}} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{H} \mathbf{y} \quad , \quad (3.9)$$

wobei $\mathbf{H} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ als “Hat”-Matrix bezeichnet wird. Die *geschätzten Residuen* sind

$$\hat{\boldsymbol{\varepsilon}} = \mathbf{y} - \hat{\mathbf{y}} = (\mathbf{I} - \mathbf{H}) \mathbf{y} \quad . \quad (3.10)$$

Es gilt: $\mathbf{X}^\top \hat{\boldsymbol{\varepsilon}} = \mathbf{X}^\top (\mathbf{I} - \mathbf{H}) \mathbf{y} = \mathbf{0}$ und $\hat{\mathbf{y}}^\top \hat{\boldsymbol{\varepsilon}} = \hat{\boldsymbol{\beta}}^\top \mathbf{X}^\top \hat{\boldsymbol{\varepsilon}} = 0$. D.h., die geschätzten Residuen sind orthogonal zu den Spalten von \mathbf{X} und zu den geschätzten Werten $\hat{\mathbf{y}}$.

Eine **geometrische Veranschaulichung** dieser Methode lässt die Hintergründe klarer erkennen: Der Erwartungswert $E(\mathbf{Y}) = \mathbf{X} \boldsymbol{\beta}$ ist eine Linearkombination der Spalten von \mathbf{X} . $\mathbf{X} \boldsymbol{\beta}$ spannt also für variierendes $\boldsymbol{\beta}$ den Raum aller möglichen Linearkombinationen auf. I.a. wird der Beobachtungsvektor \mathbf{y} nicht exakt in diesem Raum liegen, da ein zufälliger Fehler $\boldsymbol{\varepsilon}$ additiv eingeht. Sind konkrete Beobachtungen verfügbar, ist die Summe der quadrierten Residuen $S(\boldsymbol{\beta})$ so klein wie möglich, wenn $\boldsymbol{\beta}$ so gewählt wird, dass $\mathbf{X} \boldsymbol{\beta}$ jener Punkt im aufgespannten Raum ist, der am nächsten zu \mathbf{y} ist. Das ist jener Punkt, der bei orthogonaler Projektion von \mathbf{y} auf den Raum entsteht. In Abbildung 3.1 ist der „wahre“ Wert von $\mathbf{X} \boldsymbol{\beta} = E(\mathbf{Y})$ mit Punkt A bezeichnet, der beobachtete Wert \mathbf{y} von \mathbf{Y} mit Punkt B . Den Schätzwert \hat{A} für A erhält man, indem Punkt B auf die durch die Spalten von \mathbf{X} aufgespannte Ebene orthogonal projiziert wird. Der Residuenvektor $\hat{\boldsymbol{\varepsilon}}$ (in der Abbildung mit $\hat{\boldsymbol{\varepsilon}}$ bezeichnet) ist orthogonal zu dieser Ebene (vgl. Resultate vorher).

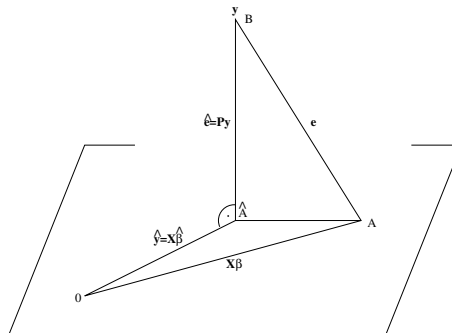
Satz 3.3.1 Für das lineare Regressionsmodell (3.2) gilt:

Der Kleinste-QS-Schätzer $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ ist unverzerrt, d.h. $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$, und $\text{Cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$.

Die Residuen $\hat{\boldsymbol{\varepsilon}}$ haben folgende Eigenschaften:

$E(\hat{\boldsymbol{\varepsilon}}) = \mathbf{0}$ und $\text{Cov}(\hat{\boldsymbol{\varepsilon}}) = \sigma^2 (\mathbf{I} - \mathbf{H})$. $\hat{\boldsymbol{\beta}}$ und $\hat{\boldsymbol{\varepsilon}}$ sind unkorreliert.

Abbildung 3.1: Geometrische Darstellung des multiplen Regressionsmodells



Beweis: siehe z.B. Johnson und Wichern (1998), S. 388.

Damit hat dieser Schätzer folgende optimale Eigenschaft:

Satz 3.3.2 (Gauss-Markov)

Beim multiplen Regressionsmodell (3.2) sei vorausgesetzt, dass die Elemente des Fehlervektors unkorreliert sind, $\text{Cov}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{I}_n$ (Homoskedastizität). Dann gilt:

- (a) $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ ist ein eindeutig bestimmter, effizienter, linearer Schätzer für $\boldsymbol{\beta}$,
- (b) $s^2 = \frac{\hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}}}{n-q-1} = \frac{1}{n-q-1} (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ ist ein unverzerrter Schätzer für die Residuenvarianz σ^2 .

Beweis: siehe z.B. Mardia et al. (1979), S. 172.

Ein effizienter Schätzer hat eine Kovarianzmatrix, die kleiner ist als die jedes anderen linearen unverzerrten Schätzers. Beim Kleinsten-QS-Schätzer spricht man auch vom *besten* linearen unverzerrten Schätzer (BLUE).

Satz 3.3.3 Gilt zu den Voraussetzungen des Gauss-Markov-Theorems (Satz 3.3.2) zusätzlich

$$\mathbf{Y} \sim N_n(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \quad , \quad (3.11)$$

so können mit der Bezeichnung $(\mathbf{X}^\top \mathbf{X})^{-1} = [(g_{ij})]$ für $i, j = 1, \dots, q$ folgende Konfidenzintervalle für die Parameter β_j ($j = 1, \dots, q$) und σ^2 ermittelt werden:

- (a) $\left[\hat{\beta}_j - t_{n-q-1; 1-\frac{\alpha}{2}} \sqrt{s^2 g_{jj}}, \hat{\beta}_j + t_{n-q-1; 1-\frac{\alpha}{2}} \sqrt{s^2 g_{jj}} \right]$
ist ein Konfidenzintervall für β_j ($j = 1, \dots, q$) mit Überdeckungswahrscheinlichkeit $1 - \alpha$.
- (b) $\left[\frac{(n-q-1)s^2}{\chi_{n-q-1; 1-\frac{\alpha}{2}}^2}, \frac{(n-q-1)s^2}{\chi_{n-q-1; \frac{\alpha}{2}}^2} \right]$
ist ein Konfidenzintervall für σ^2 mit Überdeckungswahrscheinlichkeit $1 - \alpha$.

Satz 3.3.4 (Maximum-Likelihood Schätzer)

Sei $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, wobei \mathbf{X} vollen Rang $q + 1$ habe und $\boldsymbol{\varepsilon} \sim N_n(\mathbf{0}, \sigma^2 \mathbf{I})$. Der Maximum-Likelihood (Plausible) Schätzer von $\boldsymbol{\beta}$ ist gleich dem Kleinsten-QS-Schätzer $\hat{\boldsymbol{\beta}}$. Außerdem gilt:

- (a) $\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \sim N_{q+1}(\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1})$,
- (b) $\hat{\boldsymbol{\beta}}$ und die Residuen $\hat{\boldsymbol{\varepsilon}} = \mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}}$ sind voneinander unabhängig,
- (c) $n\hat{\sigma}^2 = \hat{\boldsymbol{\varepsilon}}^\top \hat{\boldsymbol{\varepsilon}} \sim \sigma^2 \chi_{n-q-1}^2$, wobei $\hat{\sigma}^2$ der Maximum-Likelihood Schätzer von σ^2 ist.

Beweis: siehe z.B. Johnson & Wichern (1998), S. 390.

Beispiel 3.3.1 Gegeben seien die Daten mit den Prüfungsergebnissen der 88 Studenten aus Tabelle 2.1. Wir wollen eine multiple Regressionsanalyse durchführen, bei der die abhängige Größe die erste Variable Mechanik (ME) ist und die restlichen Fächer die unabhängigen Größen darstellen. Als Modell wählen wir den linearen multiplen Regressionsansatz mit einem konstanten Term.

Man erhält folgendes Resultat:

Residual Standard Error = 14.1372, Multiple R-Square = 0.3764
N = 88, F-statistic = 12.5253 on 4 and 83 df, p-value = 0

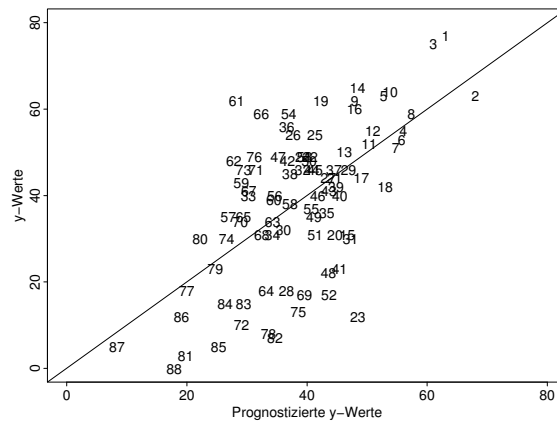
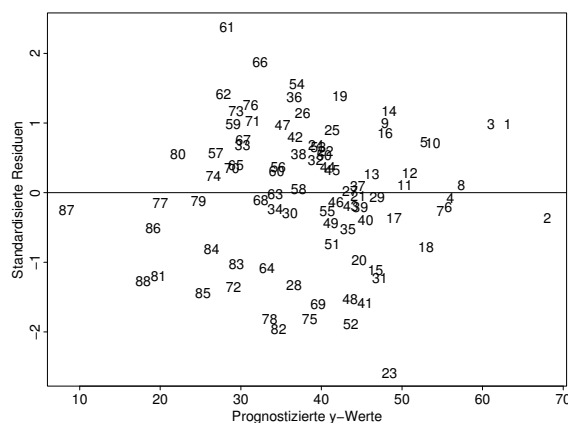
	coef	std.err	t.stat	p.value
Intercept	-12.0171	7.7542	-1.5498	0.1250
AG	0.4643	0.1461	3.1771	0.0021
LA	0.5223	0.2421	2.1572	0.0339
AN	-0.0022	0.1507	-0.0147	0.9883
ES	0.0273	0.1217	0.2241	0.8233

Die letzten beiden Fächer (AN, ES) leisten keinen Beitrag zur Erklärung der Variable ME.

Das Regressionsmodell und dessen Voraussetzungen kann mit der Residuenanalyse überprüft werden. In Abbildung 3.2 sind die geschätzten Werte $\hat{\mathbf{y}}$ den beobachteten Werten \mathbf{y} gegenübergestellt. Abbildung 3.3 zeigt $\hat{\mathbf{y}}$ gegen die standardisierten Residuen $\hat{\boldsymbol{\varepsilon}}^{norm}$ aufgetragen, wobei

$$\hat{\boldsymbol{\varepsilon}}^{norm} = \frac{\hat{\boldsymbol{\varepsilon}}}{s}$$

gilt. Zur Überprüfung auf Normalverteilung der Residuen wird in Abbildung 3.4 ein Q-Q-Plot der $\frac{i}{n+1}$ -Quantile $u_{\frac{i}{n+1}}$ der Standardnormalverteilung gegen die geordneten standardisierten Residuen $\hat{\boldsymbol{\varepsilon}}_{(1)}^{norm} \leq \hat{\boldsymbol{\varepsilon}}_{(2)}^{norm} \leq \dots \leq \hat{\boldsymbol{\varepsilon}}_{(n)}^{norm}$ dargestellt.

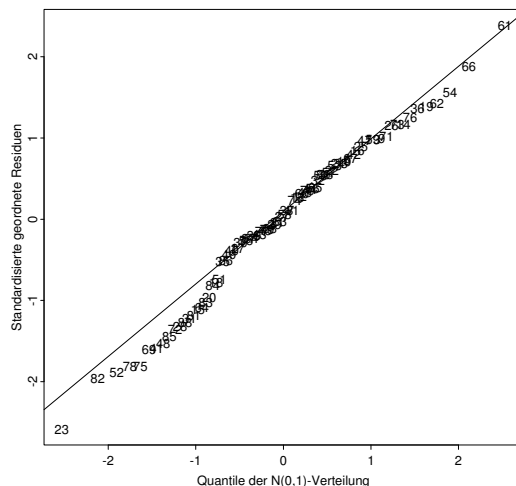
Abbildung 3.2: Residuenplot der geschätzten gegen die beobachteten y -WerteAbbildung 3.3: Residuenplot der geschätzten y -Werte gegen die standardisierten geschätzten Residuen

3.4 Multivariate lineare Regression

Haben wir bei der multiplen Regression den Fall *einer* abhängigen Größe betrachtet, so werden bei der multivariaten Regression m abhängige Größen Y_1, \dots, Y_m berücksichtigt. Somit kann für jede einzelne abhängige Größe ein Regressionsmodell für die Vorhersagevariablen x_1, \dots, x_q angesetzt werden:

$$\begin{aligned}
 Y_1 &= \beta_{01} + \beta_{11}x_1 + \dots + \beta_{q1}x_q + \varepsilon_1 \\
 &\vdots \\
 Y_j &= \beta_{0j} + \beta_{1j}x_1 + \dots + \beta_{qj}x_q + \varepsilon_j \\
 &\vdots \\
 Y_m &= \beta_{0m} + \beta_{1m}x_1 + \dots + \beta_{qm}x_q + \varepsilon_m
 \end{aligned} \tag{3.12}$$

Abbildung 3.4: Q-Q-Plot der Quantile der Standardnormalverteilung gegen die Quantile der standardisierten geschätzten Residuen



Der Fehlerterm $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_m)^\top$ hat als Voraussetzung $E(\boldsymbol{\varepsilon}) = \mathbf{0}$ und $\text{Cov}(\boldsymbol{\varepsilon}) = \boldsymbol{\Sigma}$, also die zu verschiedenen abhängigen Größen gehörigen Fehlerterme können miteinander korreliert sein.

Für eine Stichprobe von Umfang n kann nun für jede abhängige Größe $\mathbf{Y}_j = (Y_{1j}, \dots, Y_{nj})^\top$ ($j = 1, \dots, m$) das Regressionsmodell analog zum vorigen Abschnitt angeschrieben werden:

$$\begin{aligned} Y_{1j} &= \beta_{0j} + \beta_{1j}x_{11} + \dots + \beta_{qj}x_{1q} + \varepsilon_{1j} \\ Y_{2j} &= \beta_{0j} + \beta_{1j}x_{21} + \dots + \beta_{qj}x_{2q} + \varepsilon_{2j} \\ &\vdots \\ Y_{nj} &= \beta_{0j} + \beta_{1j}x_{n1} + \dots + \beta_{qj}x_{nq} + \varepsilon_{nj} \end{aligned}$$

Diese Gleichungen können mit Matrixschreibweise auch zusammengefasst werden in

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (3.13)$$

mit den $(n \times m)$ -Matrizen $\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_m)$ und $\boldsymbol{\varepsilon} = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_m)$. Aufgrund des konstanten Terms (Eins in der ersten Spalte) hat \mathbf{X} Dimension $(n \times (q + 1))$. $\boldsymbol{\beta}$ ist die $((q + 1) \times m)$ -Matrix der Regressionkoeffizienten mit den Spalten $\boldsymbol{\beta}_j$. Für den Fehlerterm gilt somit $E(\boldsymbol{\varepsilon}_j) = \mathbf{0}$ und $\text{Cov}(\boldsymbol{\varepsilon}_j, \boldsymbol{\varepsilon}_k) = \sigma_{jk}\mathbf{I}$ ($j, k = 1, \dots, m$). D.h. die m Beobachtungen des i -ten Versuches ($i = 1, \dots, n$) haben Kovarianzmatrix $\boldsymbol{\Sigma} = [(\sigma_{jk})]$, aber Beobachtungen von verschiedenen Versuchen sind unkorreliert.

Der Kleinste-QS-Schätzer bzw. der Maximum-Likelihood Schätzer für $\boldsymbol{\beta}_j$ ist

$$\hat{\boldsymbol{\beta}}_j = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}_j \quad (3.14)$$

bzw. zusammengefasst

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} . \quad (3.15)$$

Nachdem sich die Form des Regressionsschätzers zum vorigen Abschnitt nicht geändert hat, gelten auch alle Aussagen in analoger Weise, insbesondere auch der Satz von Gauss-Markov.

Literatur

- N.R. Draper and H. Smith. *Applied Regression Analysis*. Wiley & Sons, New York, 1981.
- R. Johnson and D. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, London, 4th edition, 1998.
- K.V. Mardia, J.T. Kent, and J.M. Bibby. *Multivariate Analysis*. Acad. Press, London, 1979.
- G.A.F. Seber. *Multivariate observations*. John Wiley & Sons, New York, 1984.

Kapitel 4

Auszug aus Robuster Statistik

4.1 Einleitung

Statistische Entscheidungen basieren auf Beobachtungen oder werden aufgrund spezieller Voraussetzungen (Zufälligkeit, Unabhängigkeit, Normalverteilung, usw.) getroffen. Gerade diese Voraussetzungen sind es, unter denen einerseits die klassischen Methoden funktionieren und andererseits die mathematisch einfache Handhabung ermöglicht wird. Geringfügige Abweichungen führen schon meist zu Fehlschlüssen. In der Praxis sind aber diese idealen Voraussetzungen eher selten erfüllt. Dies kann z.B. durch das Auftreten von Fehlern (Tipp- oder Übertragungsfehler) oder anderen Einflüssen zustande kommen. Auch die zugrundeliegende Verteilung wird nicht exakt mit der angenommenen übereinstimmen, sondern viel mehr in einer Umgebung dieser liegen.

Ein Ziel der robusten Statistik ist es nun Prozeduren zu finden, die gegenüber solchen Abweichungen resistent sind. Es sollen also kleine Abweichungen vom Modell auch nur geringe Auswirkungen haben.

Wir wollen uns hier auf zwei wichtige Themen konzentrieren: robuste lineare Regression und robuste Schätzung von Lokation und Kovarianz.

4.2 Robuste lineare Regression – Teil 1: Methoden

4.2.1 LS-Regression

In der Regressionsanalyse wird versucht, aufgrund gegebener Werte der Vorhersage-Variablen den genauen funktionellen Zusammenhang zur abhängigen Variable festzustellen. Man versucht also, aufgrund von bekannten Werten $1, x_{i1}, \dots, x_{iq}$ das Regressionsmodell (3.2) den beobachteten Werten y_i anzupassen. Genauer gesagt, es sollen die Regressionskoeffizienten β und die Fehlervarianz σ^2 aufgrund der vorliegenden Daten ermittelt werden.

Die Differenz

$$r_i(\beta) = y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_q x_{iq}) ,$$

die die Abweichung der Anpassung vom beobachteten Wert angibt, wird auch als *Residuum* bezeichnet. Bei der *Methode der Kleinsten Quadrate* (*Least Squares*, oder kurz LS) soll die Summe der quadrierten Residuen minimiert werden:

$$\begin{aligned} \sum_{i=1}^n r_i^2(\boldsymbol{\beta}) &= \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{i1} - \dots - \beta_q x_{iq})^2 \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \end{aligned} \quad (4.1)$$

Jener Koeffizient, der dieses Kriterium minimiert, wird mit $\hat{\boldsymbol{\beta}}_{LS}$ bezeichnet. Er ist der Kleinste-Quadratsummen-Schätzer des Regressionsparameters $\boldsymbol{\beta}$. Es ergibt sich

$$\hat{\boldsymbol{\beta}}_{LS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad . \quad (4.2)$$

Mit Hilfe von $\hat{\boldsymbol{\beta}}_{LS}$ können die *geschätzten Werte* $\hat{\mathbf{y}}_{LS}$ von \mathbf{y} berechnet werden durch

$$\hat{\mathbf{y}}_{LS} = \mathbf{X} \hat{\boldsymbol{\beta}}_{LS} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} = \mathbf{H} \mathbf{y} \quad , \quad (4.3)$$

wobei $\mathbf{H} = \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ als “Hat”-Matrix bezeichnet wird. Die *geschätzten Residuen* sind

$$\hat{\boldsymbol{\varepsilon}}_{LS} = \mathbf{y} - \hat{\mathbf{y}}_{LS} = (\mathbf{I} - \mathbf{H}) \mathbf{y} \quad . \quad (4.4)$$

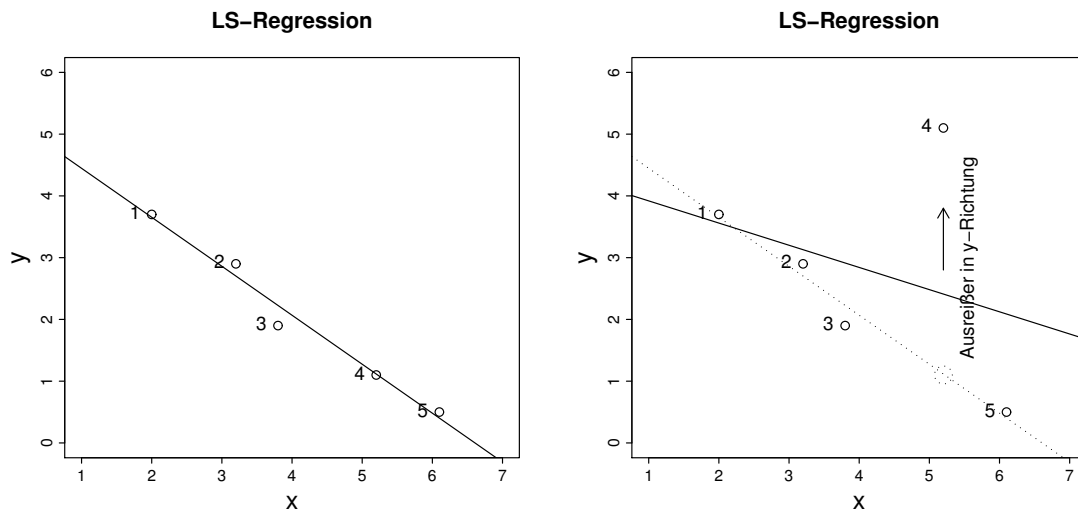
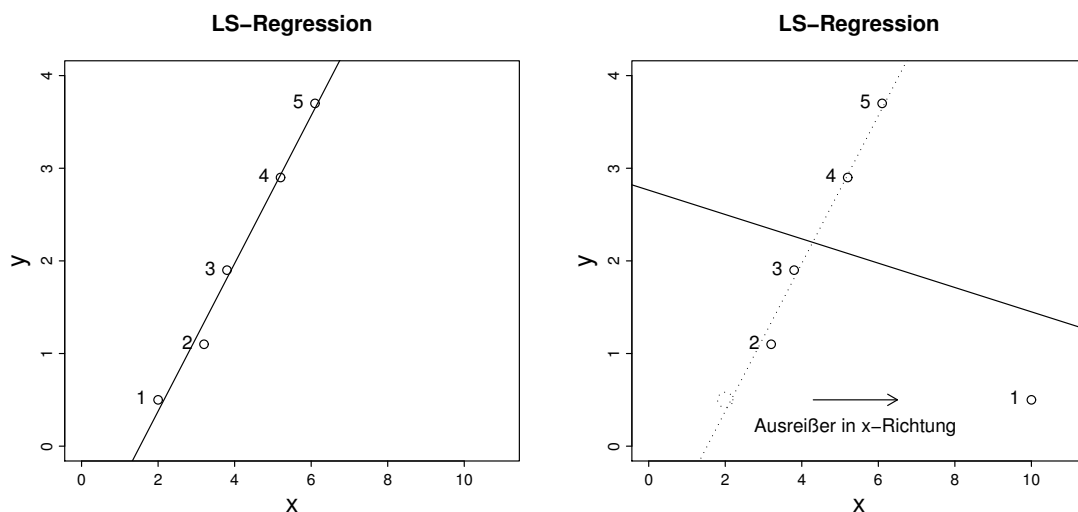
Beispiel 1: Bei der LS- oder Kleinsten-Quadrate Regression wird zur Schätzung der Koeffizienten *jeder* Wert der Residuen quadriert berücksichtigt. Das kann allerdings dazu führen, dass einzelne Werte, die nicht dem linearen Trend folgen (z.B. Ausreißer), starken Einfluss auf diese Schätzung haben. In Abbildung 4.1 wird dies illustriert. In der linken Grafik sieht man 5 Punkte, die etwa einem linearen Trend folgen. Die LS-Regressionsgerade ist eingezeichnet, die offensichtlich diese lineare Beziehung gut modelliert. Nun wird in der rechten Grafik der 4. Punkt in y -Richtung verschoben. Man erkennt, dass die LS-Regressionsparameter davon stark beeinflusst werden, weil die Gerade “nach oben gezogen” wird.

In Abbildung 4.2 wird bei einer ähnlichen Konfiguration ein Wert (Nr. 1) in x -Richtung verschoben. Die Auswirkung auf LS-Regression ist noch gravierender als vorher, weil die Regressionsgerade nicht nur abgelenkt sondern sogar “gekippt” wird. Man nennt solche x -Ausreißer *Hebelpunkte*, weil sie eine “Hebelwirkung” auf die LS-Regressionsgerade haben.

Bei der (multiplen) Regression sind *Diagnostik-Plots* üblich, um Modellfehler bzw. Ausreißer sichtbar zu machen. Der Plot der *standardisierten Residuen* gegen den Index i der Beobachtung verfolgt diesen Zweck. Die standardisierten Residuen sind definiert als

$$\frac{\hat{\varepsilon}_i(\hat{\boldsymbol{\beta}}_{LS})}{s_{LS}} \quad ,$$

also geschätzte LS-Residuen dividiert durch die geschätzte Standardabweichung der Residuen. Wenn somit obige Modellvoraussetzung $\varepsilon_i \sim N(0, \sigma^2)$ zutrifft, sollten die standardisierten Residuen etwa $N(0, 1)$ -verteilt sein. Werte, die außerhalb der in der rechten Grafik von Abbildung 4.3 eingezeichneten Grenzen ± 2.5 liegen sind somit

Abbildung 4.1: Auswirkung von y -Ausreißern bei LS-Regression.Abbildung 4.2: Auswirkung von x -Ausreißern bei LS-Regression.

potentielle Ausreißer. Abbildung 4.3 zeigt diese Plots für die LS-Regression bei den Daten mit den y -Ausreißern (links) und den x -Ausreißern (rechts). Man beachte, dass in beiden Grafiken kein Wert als Ausreißer identifiziert wird. Der Grund dafür ist die Verzerrung der Schätzung der LS-Regressionsparameter.

4.2.2 L_1 -Regression

Die Empfindlichkeit bei LS-Regression gegenüber Ausreißern ist bedingt durch das quadratische Eingehen von großen Residuen bei der zu minimierenden Funktion (4.1). Dies kann leicht vermieden werden, indem anstelle des Quadrats der Absolut-

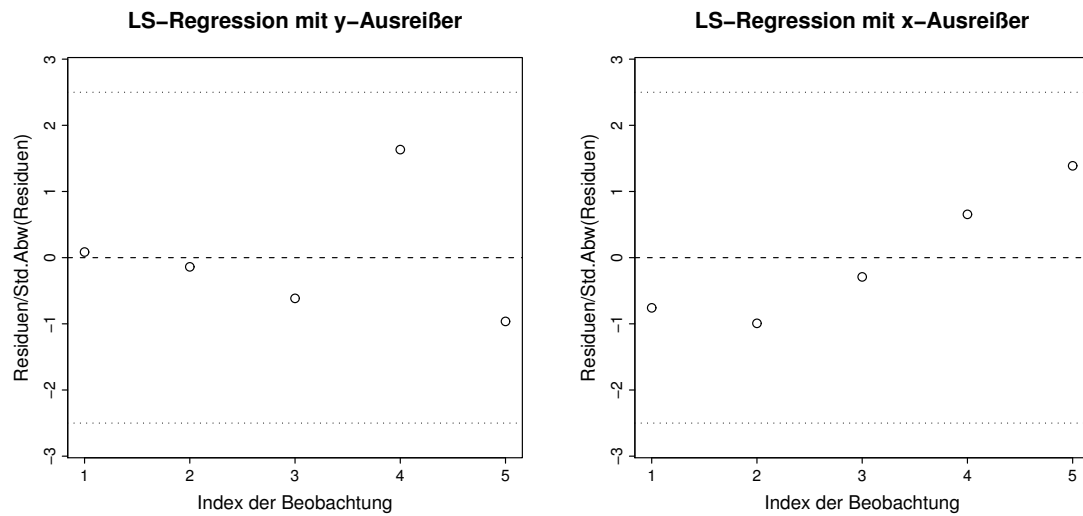


Abbildung 4.3: Plot der standardisierten Residuen auf die Daten mit y - bzw. x -Ausreißer.

betrag genommen wird. Bei L_1 -Regression wird somit

$$\sum_{i=1}^n |r_i(\beta)| \quad (4.5)$$

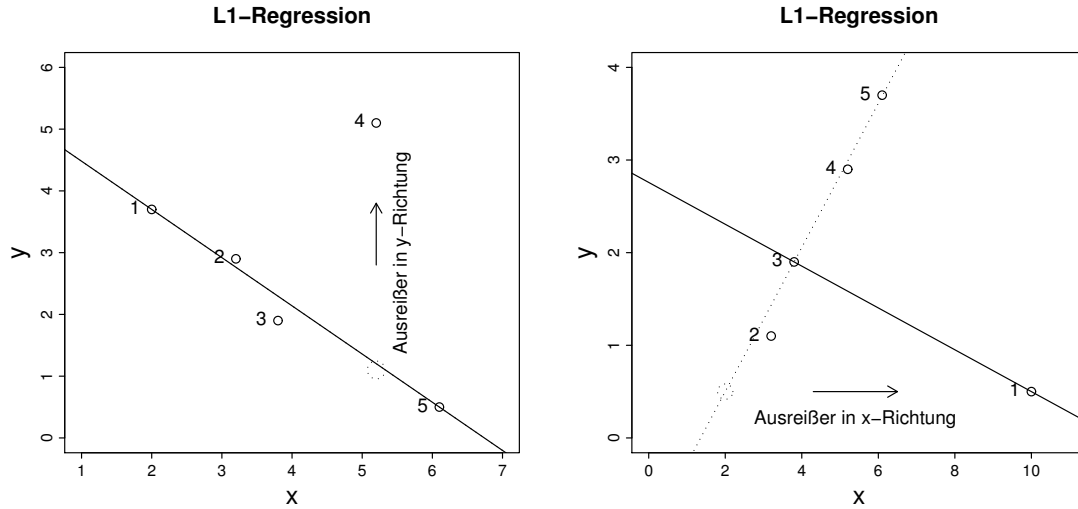
(L_1 -Norm) minimiert. Man erhält die geschätzten Regressionskoeffizienten $\hat{\beta}_{L_1}$ nicht mehr mit einer expliziten Formel, sondern muss einen Algorithmus anwenden, der die Lösung approximiert.

Beispiel 2: Wir betrachten dieselben Daten wie in Beispiel 1 für y - bzw. x -Ausreißer und schätzen die Regressionsparameter mit L_1 -Regression. Die Ergebnisse sind in Abbildung 4.4 dargestellt. Die punktierten Linien sind die jeweiligen Resultate ohne Verschiebung der Punkte, die durchgezogenen Linien beziehen sich auf die L_1 -Regressionskoeffizienten der Daten mit Ausreißern. Man erkennt, dass y -Ausreißer keinen Einfluss auf L_1 -Regression haben, während x -Ausreißer die Regressionsgerade zum Kippen bringen. Analog wie bei LS-Regression genügt also auch hier ein einziger Datenpunkt, um die Modellanpassung beliebig unsinnig zu machen. Man spricht dann auch von *Bruchpunkt* 0%, weil für ein sinnvolles Ergebnis 0% der Daten (beliebig “böswillige”) Ausreißer sein dürfen.

4.2.3 LMS-Regression und LTS-Regression

Offensichtlich hat das L_1 -Kriterium nicht zum gewünschten Erfolg geführt. Ein Grundprinzip der robusten Statistik ist es, ein Modell nicht *allen* Datenpunkten, sondern nur der Mehrheit der Daten anzupassen. Sowohl bei LS- als auch bei L_1 -Regression wurden alle Daten $i = 1, \dots, n$ berücksichtigt. Bei *Least Median of Squares* (LMS-) Regression wird die Funktion

$$\text{med}_i r_i^2(\beta) \quad (4.6)$$

Abbildung 4.4: Auswirkung von y - bzw. x -Ausreißern bei L_1 -Regression.

minimiert. Es wird also einfach die Summe in (4.1) durch den Median ersetzt. Auch hier gibt es keine explizite Lösung für die Regressionskoeffizienten $\hat{\beta}_{LMS}$, sondern nur approximative Algorithmen.

Mit Hilfe von theoretischen Überlegungen erkennt man, dass der Bruchpunkt von LMS-Regression 50% ist, d.h. bis zu 50% der Daten könnten beliebig verschoben werden ohne die Regressionskoeffizienten wesentlich zu verändern.

Eine andere sehr robuste Regressionsmethode ist *Least Trimmed Sum of Squares* (LTS-) Regression, bei der

$$\sum_{i=1}^h (r^2(\beta))_{(i)} \quad (4.7)$$

minimiert wird. Dabei sind $(r^2(\beta))_{(1)} \leq (r^2(\beta))_{(2)} \leq \dots \leq (r^2(\beta))_{(h)} \leq \dots \leq (r^2(\beta))_{(n)}$ die der Größe nach geordneten *quadratierten* Residuen. Bei einer Wahl von $h \approx n/2$ hat die Methode einen Bruchpunkt von etwa 50%, und bei größerem h sinkt der Bruchpunkt auf $(n - h)/n$.

Beispiel 3: Sowohl LMS-Regression als auch LTS-Regression (mit $3 \leq h \leq 4$) liefern stabile Resultate auf die Daten von Beispiel 2 (siehe Abbildung 4.5).

Ähnlich wie in Abbildung 4.3 bei LS-Regression kann auch hier der Plot der standardisierten Residuen angefertigt werden. Es müssen also die LMS- oder LTS-Residuen durch deren Standardabweichung dividiert werden. Die Standardabweichung σ der Residuen wird bei LMS-Regression geschätzt durch

$$\hat{\sigma}_{LMS} = c_1 \cdot \sqrt{\text{med}_i \hat{\varepsilon}_i^2(\hat{\beta}_{LMS})}$$

und bei LTS-Regression mittels

$$\hat{\sigma}_{LTS} = c_2 \cdot \sqrt{\frac{1}{h} \sum_{i=1}^h (\hat{\varepsilon}_i^2(\hat{\beta}_{LTS}))_{(i)}}.$$

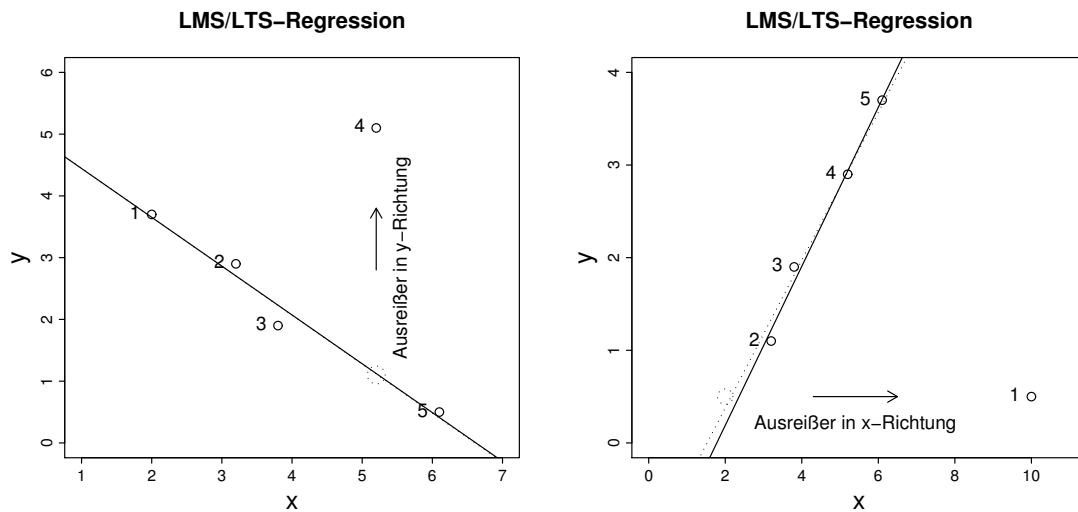


Abbildung 4.5: Auswirkung von y - bzw. x -Ausreißern bei LMS- bzw. LTS-Regression.

c_1 und c_2 sind Konstante für Konsistenz bei Normalverteilung. Man beachte, dass bei LMS- bzw. LTS-Regression nicht nur die Residuen, sondern auch die Standardabweichung der Residuen robust geschätzt werden. Somit kann der Plot der standardisierten Residuen zur Ausreißererkennung herangezogen werden. In Abbildung 4.6 sind diese Plots für die Regressionen in Abbildung 4.5 dargestellt. Man erkennt sehr deutlich, dass die Ausreißer weit außerhalb des Intervalls ± 2.5 liegen.

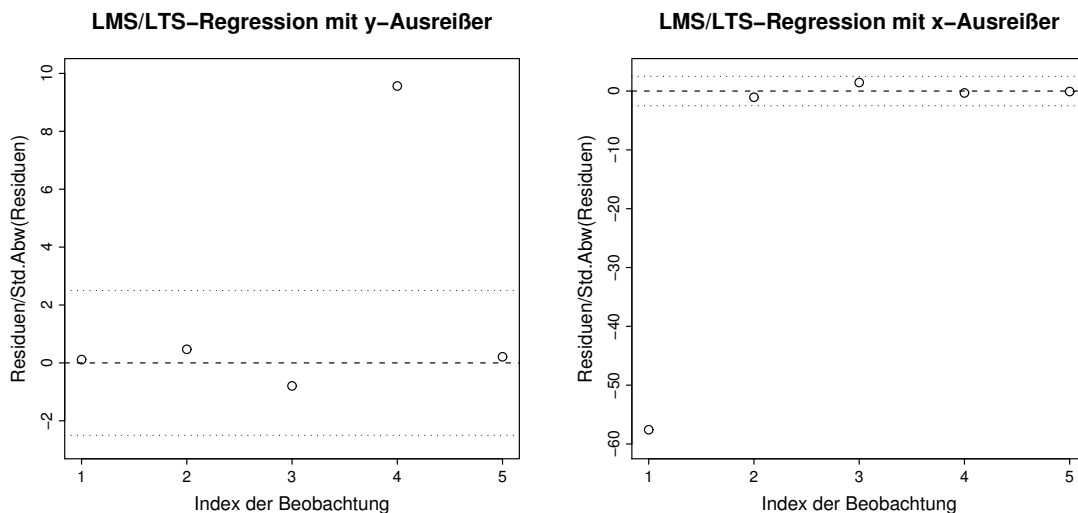


Abbildung 4.6: Plot der standardisierten Residuen auf die Daten mit y - bzw. x -Ausreißern bei LMS- bzw. LTS-Regression.

Beispiel 4: Wir betrachten die sogenannten “Stars”-Daten, die am sternförmigen Cluster CYG-OB1 von 47 Himmelskörpern in Richtung von Cygnus erhoben wurden.

Die x -Variable beschreibt das Spektrum der Himmelskörper und die y -Variable die logarithmierte Lichtintensität. Abbildung 4.7 zeigt in der linken Grafik die Daten, es wurde die LS-Regressionsgerade eingezeichnet. Man bemerkt, dass 4 Ausreißer die Schätzung so stark beeinflussen, dass das Regressionsmodell völlig unsinnig wird. In der rechten Grafik ist der Plot der standardisierten Residuen dargestellt, der keine Ausreißer erkennen lässt, weil alle Werte im Toleranzband ± 2.5 liegen.

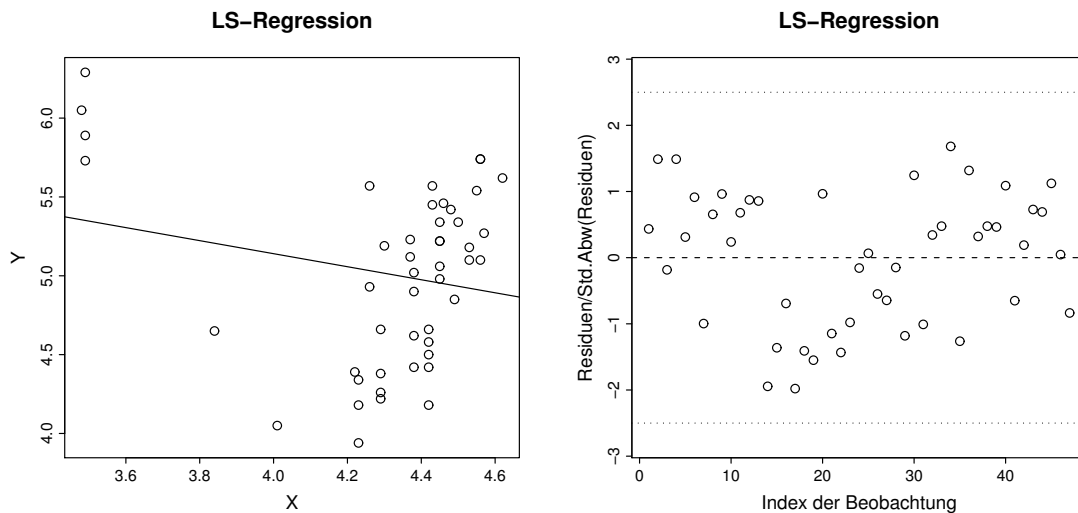


Abbildung 4.7: Stars-Daten: LS-Regression (links) und Diagnostik-Plot (rechts), der keine Ausreißer erkennen lässt.

Anders sieht dies bei LMS- oder LTS-Regression aus (Abbildung 4.8). Die Regressionsgerade folgt dem Haupttrend der Daten, aber auch der Diagnostik-Plot zeigt sehr schön die 4 extremen Ausreißer. Bei näherer Untersuchung hat sich herausgestellt, dass diese 4 Ausreißer keine Sterne, sondern Galaxien sind.

Beispiel 5: Die Daten von Hawkins, Bradu und Kass (1984) sind künstliche Daten, die zum Testen von robuster Regression konstruiert wurden. Man findet die Daten z.B. in R im Paket `rrcov` unter `hbk`. Die Daten bestehen aus einer abhängigen Größe y und drei unabhängigen Variablen x_1 , x_2 und x_3 . Bei diesem multiplen Regressionsproblem können also die Daten selbst nicht mehr wie in den vorigen Beispielen geplottet werden. Von den 75 Beobachtungen sind die ersten 14 Werte Hebelpunkte (also Ausreißer im Raum der x -Variablen), wobei aber die Werte 11-14 entlang der Regressionshyperebene liegen und somit "gute" Hebelpunkte darstellen, die die Präzision der Schätzung sogar verbessern. In Abbildung 4.9 sind die Diagnostik-Plots bezüglich LS-Regression (links) und LMS- bzw. LTS-Regression (rechts) dargestellt. Während mit LS-Regression ein Datenwert als Ausreißer erkannt wird (der aber in Richtung der Regressionshyperebene liegt), werden bei der robusten Regression genau die ersten 10 Werte als Ausreißer erkannt.

Bemerkung: Im Anschluss an LMS- oder LTS-Regression könnte zur Erhöhung der Effizienz der Schätzer eine gewichtete LS-Regression durchgeführt werden. Man

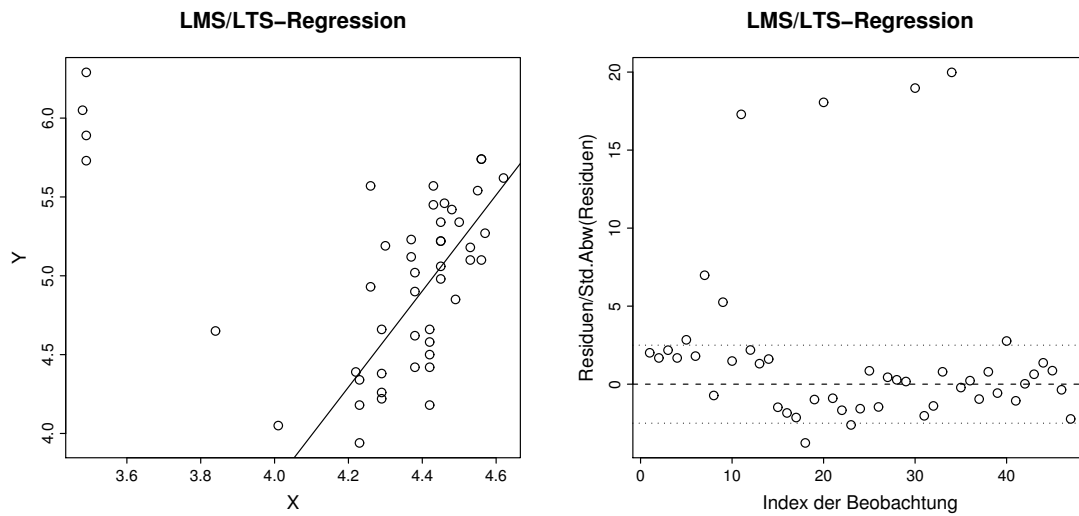


Abbildung 4.8: Stars-Daten: LTS-Regression (links) und Diagnostik-Plot (rechts), der die Ausreißer klar erkennen lässt.

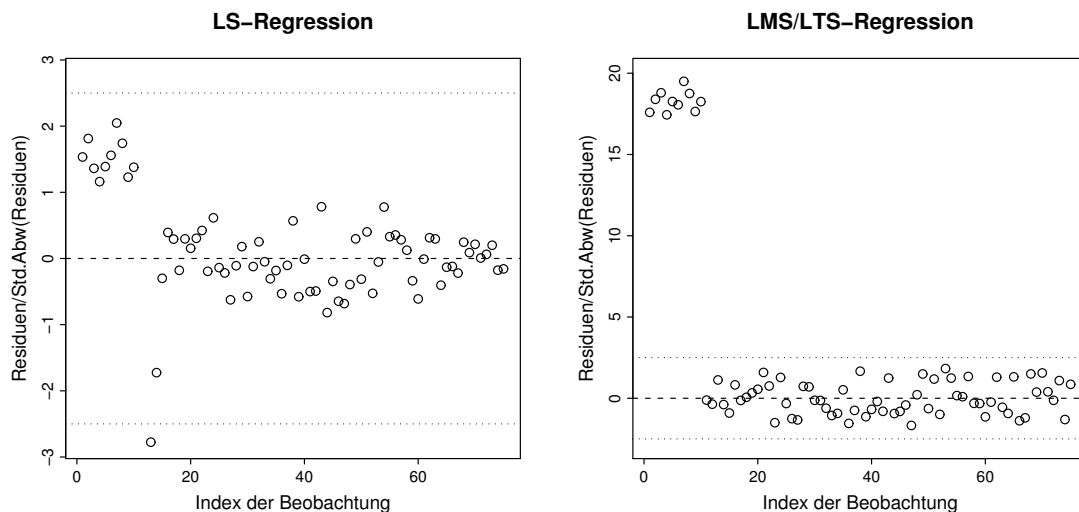


Abbildung 4.9: Daten von Hawkins, Bradu und Kass: LTS-Regression (links) und Diagnostik-Plot (rechts), der die Ausreißer klar erkennen lässt.

kann dazu jede Beobachtung mit Gewicht 0 oder 1 gewichten, je nachdem, ob diese Beobachtung im Diagnostik-Plot der robusten Regression außerhalb oder innerhalb des Toleranzbandes ± 2.5 liegt. Die gewichtete LS-Regression ist damit robust, und es können die bei LS-Regression üblichen Konfidenzintervalle bzw. Tests für die Regressionskoeffizienten verwendet werden.

4.3 Robuste Schätzung von multivariater Lokation und Kovarianz

4.3.1 Allgemeines

Die multivariate Lokation und insbesondere die Kovarianz spielen eine zentrale Rolle in der multivariaten Statistik, weil viele Methoden direkt darauf aufbauen. Bei Hauptkomponentenanalyse wird direkt die Kovarianzmatrix zur Berechnung der Hauptkomponenten herangezogen, ähnlich ist es bei der Faktorenanalyse. Bei Diskriminanzanalyse werden die Kovarianzmatrizen der einzelnen Gruppen verwendet, bei kanonischer Korrelationsanalyse gehen zusätzlich die gemischten Kovarianzen ein, usw.

Im “klassischen” Fall wird der Erwartungswert μ mit dem arithmetischen Mittel \bar{x} , und die Kovarianzmatrix Σ mit der Stichprobenkovarianzmatrix S geschätzt. Natürlich sind diese Schätzer empfindlich gegenüber Ausreißern: Das arithmetische Mittel kann bereits durch einen einzigen Ausreißer verzerrt werden, weil bei der Berechnung jeder Datenpunkt mit gleichen Gewicht eingeht (Bruchpunkt 0%). Das Gleiche gilt bei der Berechnung von S ; hier geht außerdem auch \bar{x} ein (Bruchpunkt 0%).

Es gibt viele Ansätze zur Robustifizierung der Schätzer. Wir werden uns hier auf zwei Methoden beschränken, deren Konzept sehr ähnlich ist. Wichtig ist dabei die Forderung, dass die Schätzer *affin äquivariant* sein sollen. D.h., für einen Lokationsschätzer T soll gelten

$$T(Ax_1 + b, \dots, Ax_n + b) = A \cdot T(x_1, \dots, x_n) + b$$

und für einen Kovarianz-Schätzer C

$$C(Ax_1 + b, \dots, Ax_n + b) = A \cdot C(x_1, \dots, x_n) \cdot A^T$$

für eine nicht-singuläre $(p \times p)$ -Matrix A und ein $b \in \mathbb{R}$. Das bedeutet also, dass die Schätzer auch bei Transformationen (wie z.B. Standardisierung) richtig mittransformieren.

4.3.2 MVE und MCD Schätzer

Sowohl der *Minimum Volume Ellipsoid* (MVE) Schätzer als auch der *Minimum Covariance Determinant* (MCD) Schätzer liefern robuste Schätzungen von Mittel und Kovarianz mit maximalem Bruchpunkt von 50%. Die Definition der Lokationsschätzung T und der Kovarianzschätzung C für eine Datenmatrix X ist folgendermaßen:

- MVE-Schätzer:

$T(X)$: Zentrum jenes Ellipsoids mit minimalem Volumen, das *mindestens* die Hälfte der Datenpunkte von X enthält.

$C(X)$: Gegeben durch die Form dieses Ellipsoids, multipliziert mit einem Faktor für Konsistenz bei Normalverteilung.

Formal werden die Schätzer definiert als Minimierung der Determinante von \mathbf{C} unter der Bedingung

$$\#\{i : (\mathbf{x}_i - \mathbf{T})^\top \mathbf{C}^{-1}(\mathbf{x}_i - \mathbf{T}) \leq c^2\} \geq \left\lfloor \frac{n + p + 1}{2} \right\rfloor,$$

wobei $\#$ die Kardinalität ist, und die Konstante c als $\chi_{p,0.5}^2$ gewählt wird. Die approximative Lösung wird über Ellipsoide gemacht, die durch die Kovarianzmatrix von $p + 1$ Beobachtungen bestimmt werden.

- MCD-Schätzer:

Es wird nach jenen h Punkten gesucht, deren empirische Kovarianzmatrix kleinstmögliche Determinante hat.

$\mathbf{T}(\mathbf{X})$: arithmetisches Mittel der gesuchten h Punkte.

$\mathbf{C}(\mathbf{X})$: Empirische Kovarianzmatrix der gesuchten h Punkte, multipliziert mit einem Faktor für Konsistenz bei Normalverteilung.

Wenn beim MCD-Schätzer der Parameter $h \approx n/2$ gewählt wird, erhält man den maximalen Bruchpunkt. Für größeres h sinkt der Bruchpunkt auf $(n - h)/n$.

Man kann zeigen, dass die h Punkte, die zur MCD Lösung führen, innerhalb eines Ellipsoids liegen. Dies ist aber nicht notwendigerweise das Ellipsoid mit kleinstem Volumen, womit die MCD von der MVE Lösung i.A. verschieden ist.

Beispiel 6: Wir betrachten 2-dimensionale Daten mit 18 Messungen von anorganischem und organischem Phosphor im Boden. Abbildung 4.10 zeigt die Daten, gemeinsam mit 97.5% Toleranzellipsen. Eine $(1 - \alpha)\%$ Toleranzellipse enthält im Falle von normalverteilten Daten $(1 - \alpha)\%$ der Datenpunkte. In der linken Grafik wurde die Toleranzellipse mit der Lokation und Kovarianz des MCD-Schätzers ermittelt (MVE würde ein sehr ähnliches Resultat liefern), in der rechten Grafik wurden die klassischen Schätzungen $\bar{\mathbf{x}}$ und \mathbf{S} verwendet. Datenpunkte, die außerhalb der Toleranzellipse liegen, sind potentielle Ausreißer. Man erkennt, dass nur mit der robusten Methode die offensichtlichen Ausreißer identifiziert werden.

Beispiel 6 hat gezeigt, dass mit Hilfe von robuster Lokations- und Kovarianzschätzung auch Ausreißer identifiziert werden können. Allerdings bleibt diese Methode mit der Toleranzellipse auf $p = 2$ Dimensionen beschränkt, in höherer Dimension ist das nicht mehr sinnvoll durchführbar.

Stattdessen kann aber ein äquivalenter Ansatz gewählt werden, nämlich die Berechnung der **Mahalanobis-Distanz**, deren Quadrat definiert ist als

$$\text{MD}_i^2 = (\mathbf{x}_i - \bar{\mathbf{x}})^\top \mathbf{S}^{-1}(\mathbf{x}_i - \bar{\mathbf{x}}). \quad (4.8)$$

$\bar{\mathbf{x}}$ ist das arithmetische Mittel und \mathbf{S} die empirische Kovarianzmatrix. Die Mahalanobis-Distanz wird für jeden Datenpunkt $\mathbf{x}_i \in \mathbb{R}^p$ berechnet, sie gibt an, wie weit ein Wert vom Zentrum der Punktwolke der Daten entfernt ist, relativ zur Größe der Punktwolke. Im Falle von p -variater normalverteilter Daten ist die quadrierte Mahalanobis-Distanz approximativ verteilt nach χ_p^2 . Nachdem große Werte der Mahalanobis-Distanz auf Ausreißer schließen lassen, kann ein Quantil, z.B. $\chi_{p;0.975}^2$ verwendet

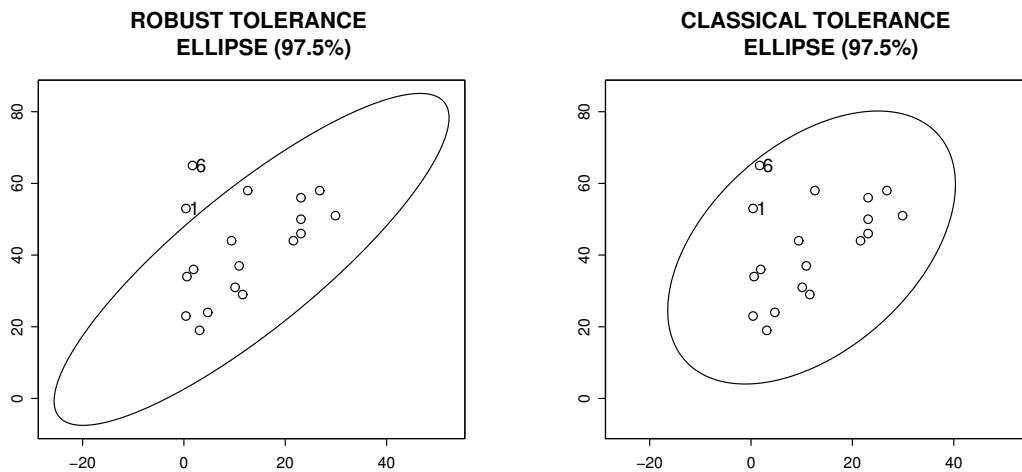


Abbildung 4.10: Phosphor-Daten: MCD-Schätzer (links) und klassische Schätzung (rechts).

werden, um Datenpunkte als Ausreißer zu deklarieren. Es kann gezeigt werden, dass diese Methode tatsächlich äquivalent ist zur Ausreißerbestimmung mittels Toleranzellipse. Der Vorteil ist aber, dass die Mahalanobis-Distanz ein 1-dimensionales Maß ist, und daher für Daten beliebiger Dimension p grafisch veranschaulicht werden kann.

Natürlich werden wir nicht die Mahalanobis-Distanz laut Definition (4.8) verwenden, weil dieses Maß, das auf klassischen Schätzern basiert, nicht robust ist. Stattdessen ersetzen wir die Schätzung von Mittel und Kovarianz durch robuste Gegenstücke, wie sie z.B. der MVE- oder MCD-Schätzer liefert. Diese robustifizierte Mahalanobis-Distanz nennen wir dann **Robuste Distanz** (RD).

Beispiel 7: Wir betrachten nochmals die Phosphor-Daten von Beispiel 6 und berechnen sowohl die Robuste Distanz als auch die Mahalanobis-Distanz (siehe Abbildung 4.11). Wie bei Beispiel 6 werden mit der Robusten Distanz die Beobachtungen 1 und 6 als Ausreißer identifiziert, da sie über dem Wert $\sqrt{\chi_{2;0.975}^2} = 2.72$ liegen. Die Mahalanobis-Distanz wird von den Ausreißern verzerrt und daher werden die Ausreißer nicht als solche angezeigt. Man spricht in diesem Fall auch von *maskierten* Ausreißern bzw. vom *Maskierungseffekt*. Dieser Effekt tritt speziell dann auf, wenn nicht einzelne sondern eine ganze Gruppe von Ausreißern auftritt.

Beide Grafiken können zu einer einzigen zusammengefasst werden, da der Index der Beobachtung als eigene Achse oft nicht so relevant ist. Der resultierende **Distance-Distance Plot** ist in Abbildung 4.12 gezeigt. Mit dem horizontal und vertikal eingezeichneten Wert $\sqrt{\chi_{2;0.975}^2} = 2.72$ werden 4 Quadranten definiert: Im Quadranten links unten befinden sich die regulären Beobachtungen, die keine Ausreißer sind. Der Quadrant rechts unten sollte leer bleiben. Rechts oben befinden sich die Ausreißer, die sowohl mit der klassischen als auch mit der robusten Methode als solche erkannt werden. Im Quadranten links oben sind schließlich die maskierten

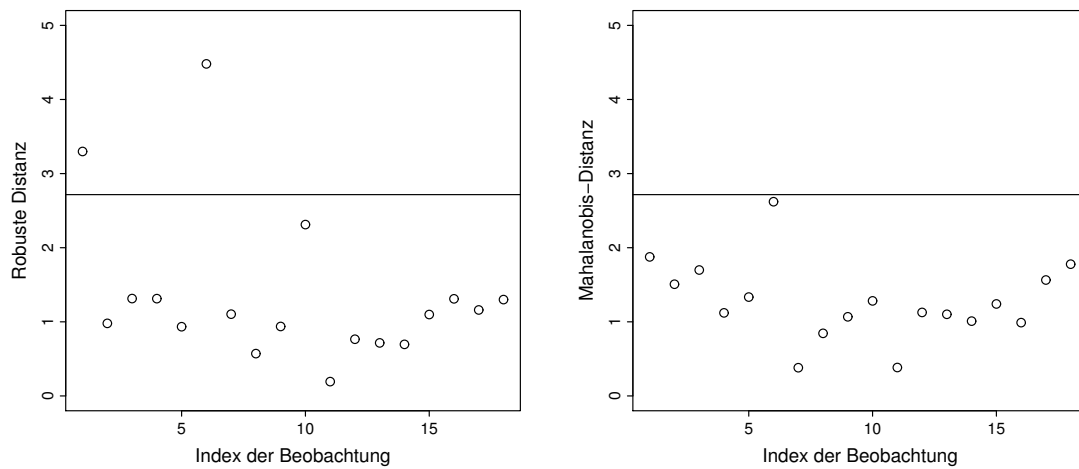


Abbildung 4.11: Phosphor-Daten: Robuste Distanzen (links) und Mahalanobis-Distanzen (rechts).

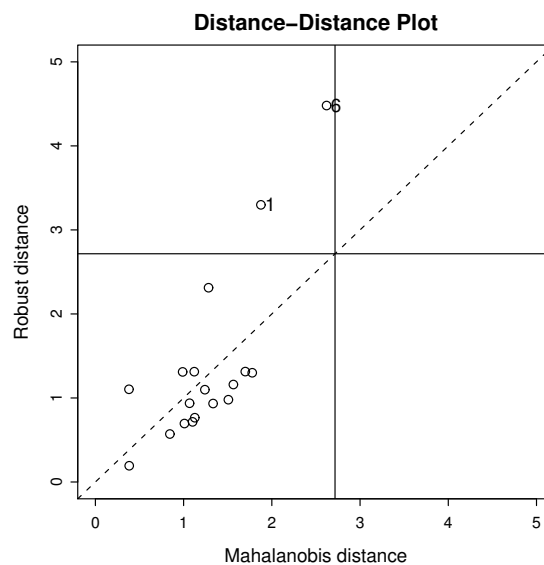


Abbildung 4.12: Phosphor-Daten: Robuste Distanzen versus Mahalanobis-Distanzen.

Ausreißer, die nur mit der Robusten Distanz erkannt werden. Wenn alle Datenpunkte auf der strichlierten Linie liegen würden, kämen die Daten genau aus einer multivariaten Normalverteilung.

Beispiel 8: Wir verwenden die sogenannten *Stackloss* Daten (in R im Package *rrcov*), die in einem kontrollierten Prozess die Pflanzenoxidation beschreiben. Es wurden an 21 aufeinanderfolgenden Tagen die Flussgeschwindigkeit der Kühlluft, die Kühlwassertemperatur und die Säurekonzentration gemessen. Der Distance-Distance Plot ist in Abbildung 4.13 links dargestellt. Mit der klassischen Mahalanobis-

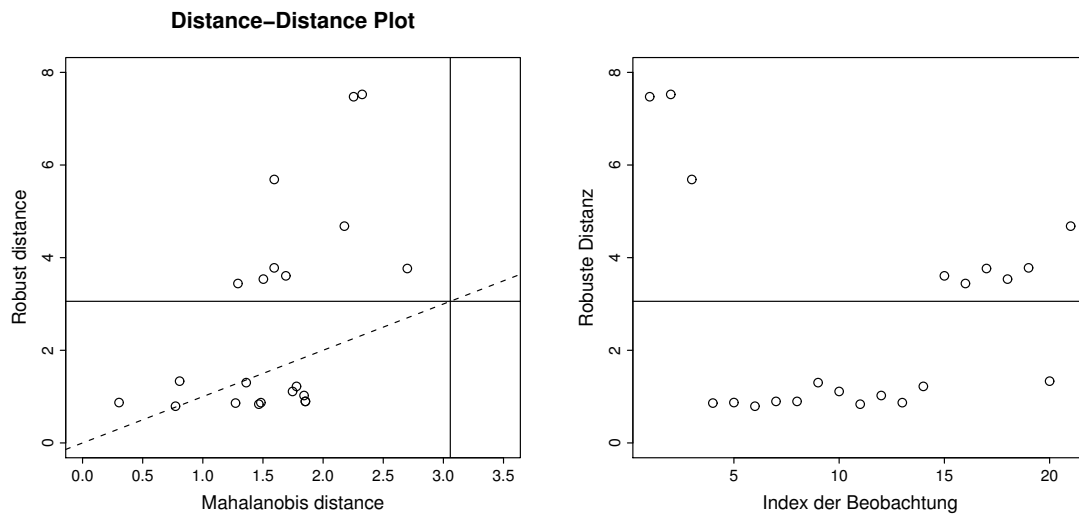


Abbildung 4.13: Stackloss Daten: Robuste Distanzen versus Mahalanobis-Distanzen.

Distanz wird kein einziger Wert als Ausreißer deklariert. Nur die Robuste Distanz zeigt eine ganze Gruppe von maskierten Ausreißern. Zur besseren Übersicht ist in der rechten Grafik die Robuste Distanz gegen den Index der Beobachtung aufgetragen. Die Ausreißer sind an den ersten 3 Tagen und in den letzten Tagen des Experiments zustande gekommen. Offenbar sind die Vorgänge hier sehr verschieden vom stabilen Zustand in der Mitte dieser Periode gewesen.

4.4 Robuste lineare Regression – Teil 2: Diagnostik

Wie bei anderen statistischen Modellanpassungen ist man auch bei der Regression daran interessiert, das Modell und die Voraussetzungen für die Anpassung zu überprüfen. Bei der Regression macht man daher Plots der Residuen, um die Voraussetzung der Normalverteilung und der Gleichheit der Residuenvarianz zu überprüfen. Von besonderem Interesse ist außerdem, ob die Daten dem Modell folgen, oder ob es etwa einzelne Beobachtungen gibt, die davon abweichen. Solche “Ausreißer”, die die klassische LS-Regression stark beeinflussen können, sind interessant, weil die Struktur dieser Daten völlig anders ist – aus welchen Gründen auch immer.

4.4.1 Hat-Matrix zur Identifikation von Hebelpunkten

Hebelpunkte (*Leverage Points*) sind Datenpunkte, die eine “Hebelwirkung” auf nicht robuste Regressionsmethoden ausüben können, also die Schätzung der Regressionsparameter stark beeinflussen können. Ein beliebtes Werkzeug für das Auffinden von Hebelpunkten sind die Diagonalelemente der Hat-Matrix aus Abbildung 4.14, definiert als $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$. Es gilt, dass $\hat{\mathbf{y}}_{LS} = \mathbf{H} \mathbf{y}$ ist, und insbesondere

ist

$$\hat{y}_{i_{LS}} = (h_{i1} \cdots \underline{h_{ii}} \cdots h_{in}) \begin{pmatrix} y_1 \\ \vdots \\ \underline{y_i} \\ \vdots \\ y_n \end{pmatrix}.$$

h_{ij} ist somit interpretierbar als der Effekt der j -ten Beobachtung y_j auf $\hat{y}_{i_{LS}}$. Speziell gibt h_{ii} den Einfluss der i -ten Beobachtung auf die eigene Vorhersage $\hat{y}_{i_{LS}}$ an.

Eigenschaften der (Diagonalelemente der) Hat-Matrix:

- $\mathbf{H}^\top = \mathbf{H}$, $\mathbf{H} = \mathbf{H}\mathbf{H}$, $\text{trace}(\mathbf{H}) = \text{trace}[\mathbf{X}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1}] = \text{trace}(\mathbf{I}_{q+1}) = q + 1$
Daraus folgt, dass das Mittel von h_{ii} gleich $\frac{q+1}{n}$ ist.
- $h_{ii} = (\mathbf{H}\mathbf{H})_{ii} = \sum_{j=1}^n h_{ij}h_{ji} = \sum_{j=1}^n h_{ij}^2 \iff h_{ii} = h_{ii}^2 + \sum_{j \neq i} h_{ij}^2$
Daraus folgt, dass $0 \leq h_{ii} \leq 1$ ist.
Ist $h_{ii} = 0$, dann ist auch $h_{ij} = 0$ für alle j , und es gibt dann keinen Beitrag von \mathbf{y} auf $\hat{y}_{i_{LS}}$.
Ist $h_{ii} = 1$, dann ist $h_{ij} = 0$ für alle $j \neq i$, und man erhält eine exakte Anpassung $\hat{y}_{i_{LS}} = y_i$.

Bei einer exakten Anpassung kann man davon ausgehen, dass es sich um einen Hebelpunkt handelt, weil ein extrem starker Einfluss auf die eigene Schätzung ausgeübt wird. I.A. sollte Vorsicht geboten werden bei großen Werten h_{ii} . Als **Faustregel** könnte eine Schranke von $h_{ii} > 2 \cdot \frac{q+1}{n}$ zur Identifizierung von Hebelpunkten genommen werden.

Beispiel 9: Wir betrachten die Daten von Hawkins, Bradu und Kass aus Beispiel 5, die so konstruiert wurden, dass die ersten 14 Werte Hebelpunkte sind. Laut unserem obigen Kriterium sind Werte dann verdächtig auf Hebelpunkte, wenn die Diagonalelemente der Hat-Matrix den Wert $2 \cdot \frac{q+1}{n} = 2 \cdot \frac{4}{75} = 0.107$ überschreiten. In Abbildung 4.14 sind die Werte h_{ii} gegen den Index $i = 1, \dots, 75$ aufgetragen. Die Schranke 0.107 ist als strichlierte Linie eingezeichnet. Man erkennt, dass nur die Werte 12, 13 und 14 als potentielle Hebelpunkte erkannt werden, die aber “gute” Hebelpunkte darstellen (d.h. deren \mathbf{x}_i Werte sind Ausreißer, aber die entsprechenden y_i Werte folgen dem Modell).

4.4.2 Robuste Distanz zur Identifikation von Hebelpunkten

In obigem Beispiel wurden durch die Verwendung der Hat-Matrix die “schlechten” Hebelpunkte maskiert, da das Kriterium offenbar nicht robust ist. Der Grund dafür ist leicht einzusehen. Man kann nämlich zeigen, dass bei Regression mit konstantem Term gilt:

$$h_{ii} = \frac{1}{n-1} \text{MD}_i^2 + \frac{1}{n} \quad (4.9)$$

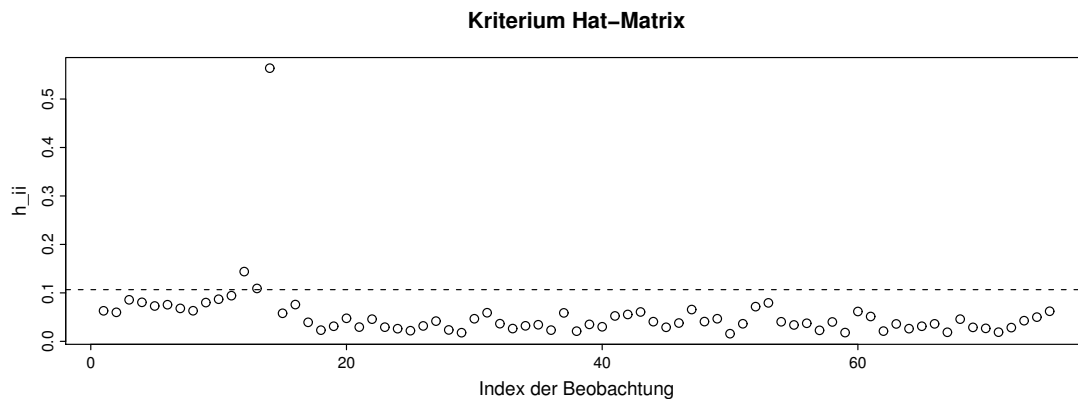


Abbildung 4.14: Daten von Hawkins-Bradru-Kass: Werte über der strichlierten Linie sind laut Hat-Matrix Kriterium potentielle Hebelpunkte.

Die Diagonalelemente der Hat-Matrix sind somit direkt proportional zur quadrierten Mahalanobis-Distanz. Von den Eigenschaften der Mahalanobis-Distanz folgt somit, dass die Hat-Matrix nur dann nicht durch Ausreißer maskiert wird, wenn es einen einzelnen Ausreißer in x -Richtung gibt, wenn aber keine Gruppe von x -Ausreißern vorliegt.

Bemerkung: Die Werte h_{ii} ziehen nur die x -Werte, nicht aber die y -Werte in Betracht. Es wird also nicht zwischen guten und schlechten Hebelpunkten unterschieden.

Beispiel 10: Um die Beziehung (4.9) plausibel zu machen, werden für die Hawkins-Bradru-Kass Daten die Mahalanobis-Distanzen berechnet und analog zu Abbildung 4.14 in Abbildung 4.15 visualisiert. Man erkennt das gleiche Bild wie bei Verwendung der Hat-Matrix, und erhält somit die gleiche verzerrte Aussage bezüglich Hebelpunkte.

Um dem Maskierungseffekt bei der Hat-Matrix bzw. bei der Mahalanobis-Distanz zu entgehen, muss das Kriterium robustifiziert werden. Es liegt aber nun auf der Hand, wie dies gemacht wird, nämlich einfach mit der Robusten Distanz.

Beispiel 11: Wir berechnen für die Hawkins-Bradru-Kass Daten die Robusten Distanzen RD_i , indem Lokation und Kovarianz durch den MCD-Schätzer berechnet werden. Das Resultat ist in Abbildung 4.16 dargestellt. Alle Hebelpunkte 1-14 werden nun richtig durch hohe Werte RD_i angezeigt.

4.4.3 Diagnostik bei Regression

In den vorigen Abschnitten haben wir nun verschiedene Fälle von Ausreißern und Hebelpunkten bei Regression kennengelernt. Wir können daher 4 Typen von Datenwerten unterscheiden:

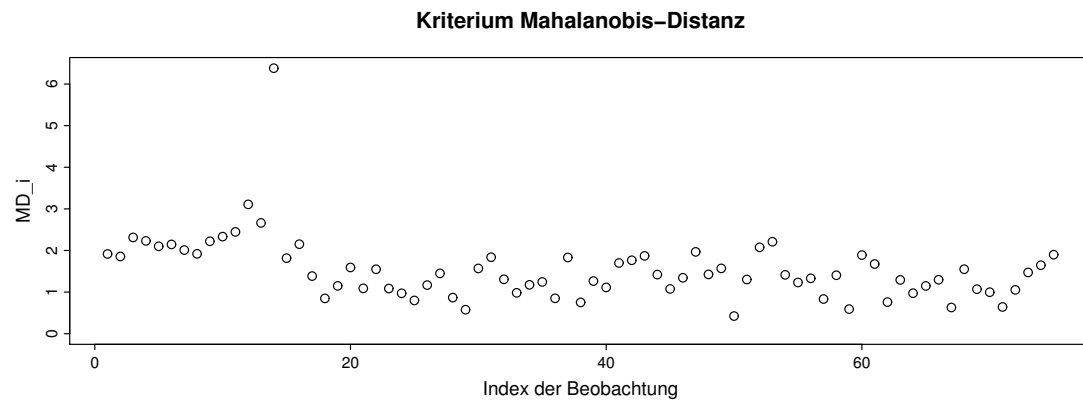


Abbildung 4.15: Daten von Hawkins-Bradru-Kass: Mahalanobis-Distanzen

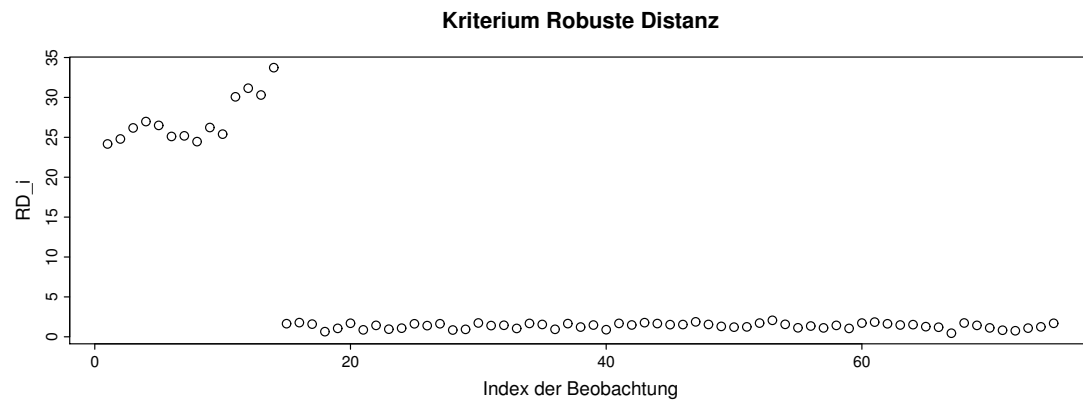


Abbildung 4.16: Daten von Hawkins-Bradru-Kass: Robuste Distanzen

- Reguläre Beobachtungen: \mathbf{x}_i im Datenbereich und y_i zum Modell passend
- Vertikale Ausreißer: \mathbf{x}_i im Datenbereich aber y_i nicht zum Modell passend
- Gute Hebelpunkte: Ausreißer bei \mathbf{x}_i aber y_i zum Modell passend
- Schlechte Hebelpunkte: Ausreißer bei \mathbf{x}_i und y_i nicht zum Modell passend

Diese 4 Typen werden im Fall von einfacher linearer Regression in Abbildung 4.17 illustriert. Im Allgemeinen bringen gute Hebelpunkte einen Vorteil, weil sie in Richtung der Regressionsgeraden liegen und damit die Schätzung der Regressionsparameter sogar noch präziser machen. Schlechte Hebelpunkte können starken Einfluss auf die Schätzung der Parameter haben, weil sie die Regressionsgerade zum Kippen bringen können.

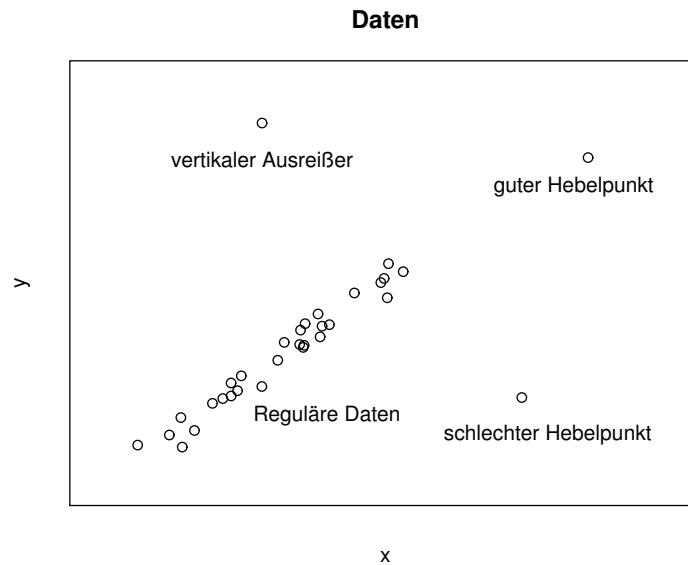


Abbildung 4.17: 4 Typen von Beobachtungen bei einfacher linearer Regression.

Wir haben in den vorhergehenden Abschnitten den Diagnostik-Plot kennengelernt, bei dem die standardisierten LMS- oder LTS-Residuen $\varepsilon_i/\hat{\sigma}$ gegen den Index der Beobachtung aufgetragen werden. Wir haben aber auch den Distance-Distance Plot kennengelernt, bei dem die Robuste Distanz der Mahalanobis-Distanz gegenübergestellt wird. Zur Unterscheidung der erwähnten 4 Typen von Beobachtungen können beide Plots kombiniert werden.

Regression Diagnostic Plot: Ausreißer sind bei Regression gekennzeichnet durch einen hohen Wert der standardisierten robusten Residuen (außerhalb der Schranken ± 2.5). Hebelpunkte sind x -Ausreißer, die somit große Robuste Distanz haben. Der *Regression Diagnostic Plot* macht eine Gegenüberstellung der standardisierten robusten Residuen und der Robusten Distanz.

Abbildung 4.18 zeigt den Regression Diagnostic Plot von den Daten aus Abbildung 4.17. Die regulären Beobachtungen, aber auch die guten Hebelpunkte, findet man im Band ± 2.5 der standardisierten robusten Residuen. Die Unterscheidung, ob eine Beobachtung Hebelpunkt oder kein Hebelpunkt ist, wird aufgrund des Wertes $\sqrt{\chi_{q;0.975}^2}$ bei der Robusten Distanz gemacht.

Beispiel 12: Zu den Phosphor-Daten von Beispiel 6 wurden auch y -Werte ermittelt, nämlich der Phosphor-Gehalt in Maispflanzen, die in den Böden mit den entsprechenden Werten (x_{i1}, x_{i2}) von Abbildung 4.10 gewachsen sind. Der Regression Diagnostik Plot ist in Abbildung 4.19 dargestellt. In der linken Grafik sieht man den Plot mit den robusten Schätzungen, in der rechten Grafik ist zum Vergleich der Plot mit den klassischen Schätzungen gemacht worden. Analog zu Abbildung 4.11 bzw. 4.12 werden x -Ausreißer aufgrund eines Wertes der Robusten Distanz größer als $\sqrt{\chi_{2;0.975}^2} = 2.72$ erkannt. Die Beobachtungen 1 und 6 sind somit Hebelpunkte.

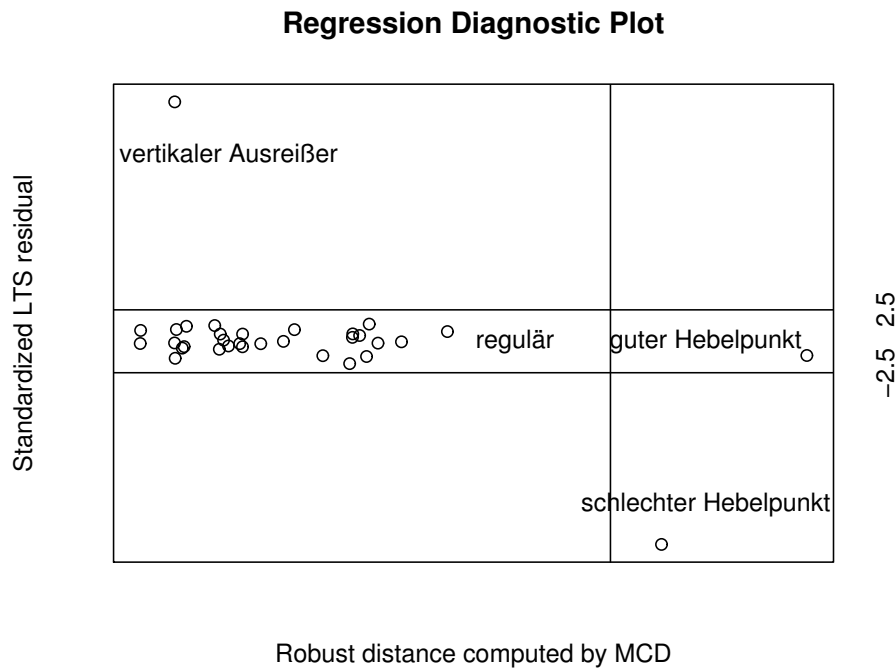


Abbildung 4.18: Identifikation der 4 Typen von Beobachtungen im Regression Diagnostic Plot.

Sie sind allerdings gute Hebelpunkte, weil sie innerhalb des Bandes ± 2.5 liegen. Der Wert 17 wird als vertikaler Ausreißer identifiziert. Auch im Plot mit den klassischen Schätzungen wird Wert 17 als vertikaler Ausreißer erkannt, es werden aber keine Hebelpunkte gefunden, diese werden maskiert in diesem “Diagnostik”-Plot.

Beispiel 13: Auch zu den *Stackloss*-Daten aus Beispiel 8 gibt es y -Werte. Der Regression Diagnostic Plot ist in Abbildung 4.20 gezeigt, wiederum mit den robusten (links) und den klassischen (rechts) Schätzungen. Wir erkennen natürlich die auch in Abbildung 4.13 mit der Schranke $\sqrt{\chi_{3,0.975}^2} = 3.06$ identifizierten x -Ausreißer als Hebelpunkte. Die Werte 1, 3 und 21 sind dabei schlechte Hebelpunkte. Ein Wert (Nr. 4) wird als vertikaler Ausreißer erkannt. Man bemerke, dass im “klassischen Plot” rechts alle Werte als reguläre Beobachtungen eingestuft werden.

Literatur

- N.R. Draper and H. Smith. *Applied Regression Analysis*. Wiley & Sons, New York, 1981.
- D.M. Hawkins, D. Bradu, and G.V. Kass. Location of several outliers in multiple regression data using elemental sets. *Technometrics*, 26, 197-208, 1984.
- R. Johnson and D. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, London, 4th edition, 1998.

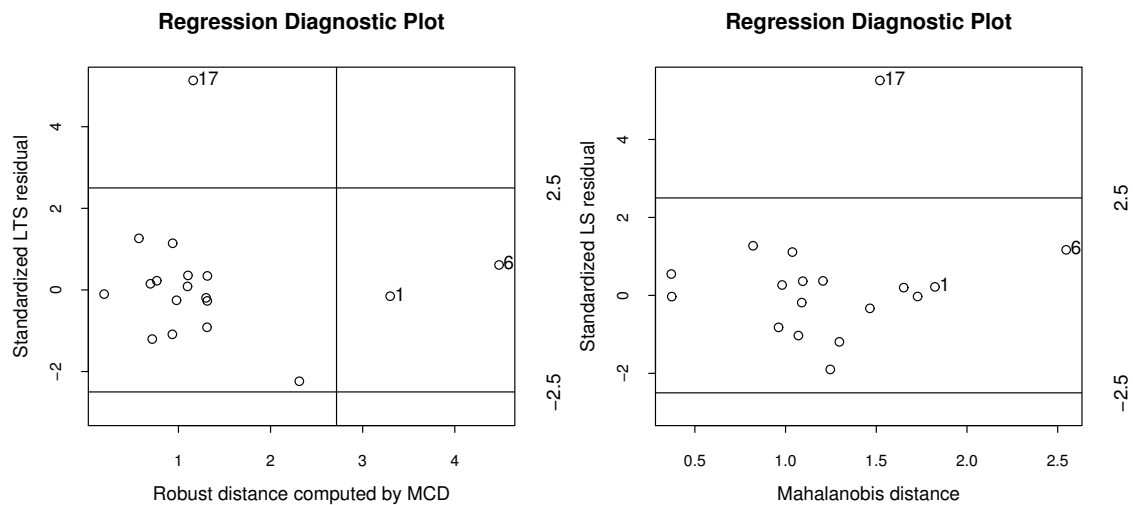


Abbildung 4.19: Regression Diagnostik Plot der Phosphor-Daten: links mit den robusten Schätzungen, rechts mit den klassischen Schätzungen.

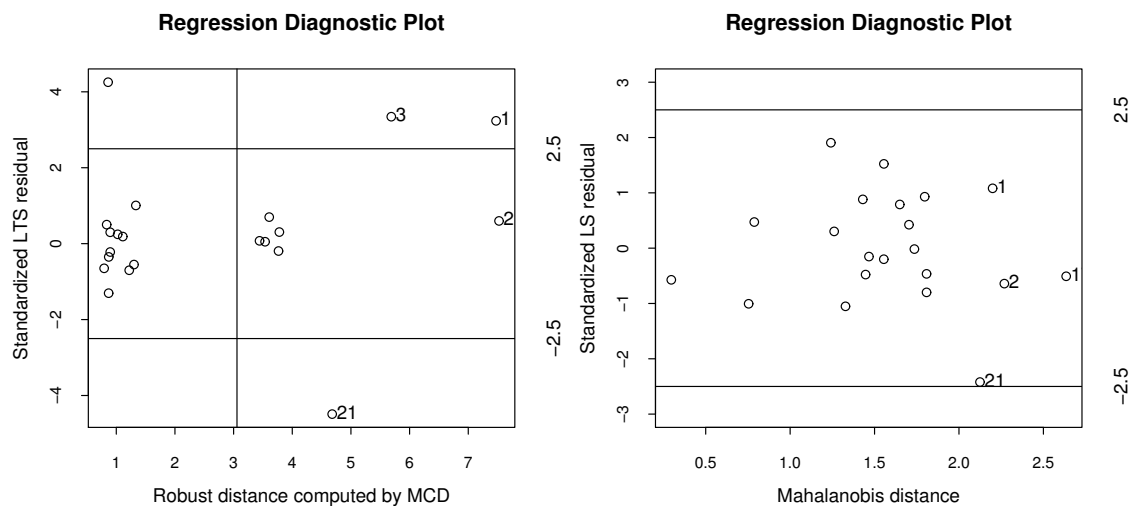


Abbildung 4.20: Regression Diagnostik Plot der Stackloss-Daten: links mit den robusten Schätzungen, rechts mit den klassischen Schätzungen.

- K.V. Mardia, J.T. Kent, and J.M. Bibby. *Multivariate Analysis*. Acad. Press, London, 1979.
- P.J. Rousseeuw. Least median of squares regression. *Journal of the American Statistical Association*, 79(388), 871-880, 1984.
- P.J. Rousseeuw. Multivariate estimation with high breakdown point. In: W. Grossmann, G. Pflug, I. Vincze, and W. Wertz (eds.), *Mathematical Statistics and Applications*, Volume B. Akadémiai Kiadó, Budapest, pp. 283-297, 1985.
- P.J. Rousseeuw and A.M. Leroy. *Robust Regression and Outlier Detection*. Wiley, New York, 1987.
- G.A.F. Seber. *Multivariate Observations*. John Wiley & Sons, New York, 1984.

Kapitel 5

Hauptkomponentenanalyse

5.1 Einleitung

Die Hauptkomponentenanalyse (Principal Component Analysis, *PCA*) hat ihren Ursprung bei Karl Pearson (1901) und wurde von Harrold Hotelling (1933) entwickelt.

Ziel dieser Methode ist es, komplizierte Beziehungen in beobachteten Daten auf eine einfache Form zu reduzieren. Die Daten werden durch Linearkombinationen von gewissen Komponenten so dargestellt, dass möglichst wenig Information verloren geht. Die Dimension wird somit reduziert auf die Anzahl dieser Komponenten.

Die Hauptkomponentenanalyse dient auch zur Ausreißererkennung. Dazu werden alle verschiedenen Paare von Komponenten einander gegenübergestellt und diese Projektionen untersucht.

Die Hauptkomponentenanalyse gehört zu den wichtigsten Methoden in der multivariaten Statistik. Sie ist einerseits relativ einfach in der Handhabung und Interpretation und bildet auf der anderen Seite den Grundstein für viele multivariate Methoden.

5.2 Bestimmung von Hauptkomponenten in der Population

Sei $\mathbf{x} = (x_1, \dots, x_p)^\top$ ein p -dimensionaler Zufallsvektor mit Mittelvektor $\boldsymbol{\mu}$ und Kovarianzmatrix

$$\boldsymbol{\Sigma} = E\left[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top\right]. \quad (5.1)$$

Sei weiters $\boldsymbol{\Gamma} = (\boldsymbol{\gamma}_1, \dots, \boldsymbol{\gamma}_p)$ eine $(p \times p)$ -Matrix mit festen Koeffizienten mit der Bedingung, dass $\boldsymbol{\gamma}_i$ Einheitsvektoren sind, d.h. $\boldsymbol{\gamma}_i^\top \boldsymbol{\gamma}_i = 1$ für $i = 1, \dots, p$. Außerdem sei $\boldsymbol{\Gamma}$ orthogonal, also $\boldsymbol{\Gamma}^\top = \boldsymbol{\Gamma}^{-1}$. Dann betrachten wir die Lineartransformation

$$\mathbf{z} = \boldsymbol{\Gamma}^\top (\mathbf{x} - \boldsymbol{\mu}) \quad (5.2)$$

oder, ausgedrückt in Komponenten,

$$z_i = \boldsymbol{\gamma}_i^\top (\mathbf{x} - \boldsymbol{\mu}) \quad \text{für} \quad i = 1, \dots, p. \quad (5.3)$$

Das Resultat der obigen Transformation ist ein neuer Zufallsvektor \mathbf{z} der Dimension p . Die Varianz von z_i ($i = 1, \dots, p$) ist

$$\text{Var}(z_i) = E\left[\boldsymbol{\gamma}_i^\top (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\gamma}_i\right] = \boldsymbol{\gamma}_i^\top \boldsymbol{\Sigma} \boldsymbol{\gamma}_i. \quad (5.4)$$

Wir möchten eine Transformation erhalten, sodass die Varianz jeweils maximiert wird mit der Nebenbedingung von Einheitsvektoren $\boldsymbol{\gamma}_i$. Mathematisch kann man diese Forderung als Maximierung eines Lagrangeschen Ausdrucks mit entsprechender Randbedingung formulieren.

1. Hauptkomponente:

Für $i = 1$ möchte man also eine Komponente z_1 erhalten, sodass $\text{Var}(z_1)$ maximiert wird, und $\boldsymbol{\gamma}_1^\top \boldsymbol{\gamma}_1 = 1$ ist. Das entsprechende Lagrange-Problem lautet somit:

$$\phi_1 = \boldsymbol{\gamma}_1^\top \boldsymbol{\Sigma} \boldsymbol{\gamma}_1 - a_1(\boldsymbol{\gamma}_1^\top \boldsymbol{\gamma}_1 - 1) \quad (5.5)$$

Die partiellen Ableitungen nach den Unbekannten $\boldsymbol{\gamma}_1$ werden Null gesetzt. Man erhält

$$\frac{\partial \phi_1}{\partial \boldsymbol{\gamma}_1} = 2\boldsymbol{\Sigma} \boldsymbol{\gamma}_1 - 2a_1 \boldsymbol{\gamma}_1 = \mathbf{0} \quad (5.6)$$

bzw.

$$\boldsymbol{\Sigma} \boldsymbol{\gamma}_1 = a_1 \boldsymbol{\gamma}_1. \quad (5.7)$$

Die Lösung des Maximierungsproblems ergibt somit, dass $\boldsymbol{\gamma}_1$ Eigenvektor von $\boldsymbol{\Sigma}$ zum Eigenwert a_1 ist. Allerdings hat $\boldsymbol{\Sigma}$ insgesamt p Eigenvektoren. Welchen sollte man nun nehmen?

Mit (5.7) sieht man rasch, dass

$$\text{Var}(z_1) = \boldsymbol{\gamma}_1^\top (\boldsymbol{\Sigma} \boldsymbol{\gamma}_1) = \boldsymbol{\gamma}_1^\top (a_1 \boldsymbol{\gamma}_1) = a_1 \boldsymbol{\gamma}_1^\top \boldsymbol{\gamma}_1 = a_1$$

ist, und nachdem diese Varianz maximiert werden sollte, nimmt man jenen Eigenvektor $\boldsymbol{\gamma}_1$ zum größten Eigenwert a_1 . Die Komponente z_1 heißt dann *erste Hauptkomponente*, und $\boldsymbol{\gamma}_1$ ist die Richtung dieser Komponente.

2. Hauptkomponente:

Im nächsten Schritt, für $i = 2$, möchte man wiederum $\text{Var}(z_2)$ maximieren, unter der Bedingung $\boldsymbol{\gamma}_2^\top \boldsymbol{\gamma}_2 = 1$. Zusätzlich fordert man, dass z_1 und z_2 unkorreliert sind. Letztere Bedingung bedeutet:

$$\begin{aligned} \text{Cov}(z_1, z_2) &= \text{Cov}(\boldsymbol{\gamma}_1^\top (\mathbf{x} - \boldsymbol{\mu}), \boldsymbol{\gamma}_2^\top (\mathbf{x} - \boldsymbol{\mu})) = \text{Cov}(\boldsymbol{\gamma}_1^\top \mathbf{x}, \boldsymbol{\gamma}_2^\top \mathbf{x}) = \boldsymbol{\gamma}_1^\top \boldsymbol{\Sigma} \boldsymbol{\gamma}_2 = \\ &= \boldsymbol{\gamma}_2^\top \boldsymbol{\Sigma} \boldsymbol{\gamma}_1 = \boldsymbol{\gamma}_2^\top a_1 \boldsymbol{\gamma}_1 = a_1 \boldsymbol{\gamma}_2^\top \boldsymbol{\gamma}_1 = 0 \end{aligned}$$

Nachdem $a_1 \neq 0$ ist, ist Unkorreliertheit gleichbedeutend mit Orthogonalität von $\boldsymbol{\gamma}_1$ und $\boldsymbol{\gamma}_2$.

Man kann nun wieder das Lagrange Problem formulieren,

$$\phi_2 = \boldsymbol{\gamma}_2^\top \boldsymbol{\Sigma} \boldsymbol{\gamma}_2 - a_2(\boldsymbol{\gamma}_2^\top \boldsymbol{\gamma}_2 - 1) - b \boldsymbol{\gamma}_2^\top \boldsymbol{\gamma}_1 \quad (5.8)$$

mit den Lagrange Koeffizienten a_2 und b . Die partiellen Ableitungen nach den Unbekannten γ_2 werden Null gesetzt. Man erhält

$$\frac{\partial \phi_2}{\partial \gamma_2} = 2\Sigma\gamma_2 - 2a_2\gamma_2 - b\gamma_1 = \mathbf{0} \quad (5.9)$$

Multiplikation von links mit γ_1^\top ergibt

$$2\gamma_1^\top \Sigma \gamma_2 - 2a_2 \gamma_1^\top \gamma_2 - b \gamma_1^\top \gamma_1 = 0 - 0 - b \cdot 1 = 0 ,$$

und damit muss $b = 0$ sein. Somit reduziert sich (5.9) auf

$$\Sigma \gamma_2 = a_2 \gamma_2 ,$$

und somit ist γ_2 Eigenvektor von Σ zum nächst größten Eigenwert a_2 , und z_2 wird als *zweite Hauptkomponente* bezeichnet.

Weitere Hauptkomponenten:

Die k -te Hauptkomponente ($2 < k \leq p$) wird analog zu oben definiert, nämlich über Maximierung von $\text{Var}(z_k)$ und der Bedingung $\gamma_k^\top \gamma_k = 1$, sowie Unkorreliertkheit zu allen vorigen Hauptkomponenten, also $\text{Cov}(z_k, z_j) = 0$ für $k > j$, gleichbedeutend mit Orthogonalität $\gamma_k^\top \gamma_j = 0$.

Wenig überraschend ist die Lösung, dass γ_k Eigenvektor von Σ zum k -ten größten Eigenwert a_k ist.

Sammelt man nun alle gefundenen Eigenvektoren als Spalten in der Matrix Γ , also $\Gamma = (\gamma_1, \dots, \gamma_p)$, und die entsprechenden Eigenwerte a_1, \dots, a_p der Größe nach sortiert ($a_1 \geq a_2 \geq \dots \geq a_p \geq 0$) in der Diagonale der Matrix \mathbf{A} , also $\mathbf{A} = \text{Diag}(a_1, \dots, a_p)$, dann kann die Hauptkomponenten-Lösung in Matrixschreibweise ausgedrückt werden als

$$\Sigma \Gamma = \Gamma \mathbf{A} \quad (5.10)$$

oder auch $\Sigma = \Gamma \mathbf{A} \Gamma^\top$, siehe Spektralzerlegungssatz.

Mit diesen Definitionen ist die Lineartransformation

$$\mathbf{z} = \Gamma^\top (\mathbf{x} - \boldsymbol{\mu}) \quad (5.11)$$

bekannt unter dem Namen *Hauptkomponententransformation*, das i -te Element z_i des Vektors \mathbf{z} nennt man *i -te Hauptkomponente (HK)*.

Der Erwartungswert der HK ist Null, da

$$E(\mathbf{z}) = \Gamma^\top [E(\mathbf{x} - \boldsymbol{\mu})] = \mathbf{0}, \quad (5.12)$$

und die Kovarianz ist

$$\text{Cov}(\mathbf{z}) = \Gamma^\top \text{Cov}(\mathbf{x} - \boldsymbol{\mu}) \Gamma = \Gamma^\top \Sigma \Gamma = \mathbf{A} . \quad (5.13)$$

Die Varianz der i -ten HK entspricht also a_i , dem i -ten Eigenwert von Σ . Verschiedene HK sind unkorreliert, weil \mathbf{A} eine Diagonalmatrix ist. Die Gesamtvarianz der

HK ist die Summe aller Eigenwerte. Dies ist natürlich auch gleichzeitig die Gesamtvarianz oder Gesamtvariation von \mathbf{x} .

Die Matrix $\mathbf{\Gamma}$, die die Beziehung zwischen \mathbf{x} und \mathbf{z} herstellt, wird auch *Ladungsmatrix* (engl. loadings matrix) genannt. Deren Elemente γ_{ij} widerspiegeln den Einfluss von x_i auf z_j .

Ein interessantes Maß für den Zusammenhang zwischen \mathbf{x} und \mathbf{z} ist auch die Korrelation. Speziell das Quadrat dieser Korrelation stellt als Bestimmtheitsmaß die erklärte Variation von z_j durch x_i dar. Die Kovarianz zwischen den Variablen und den HK ist

$$\begin{aligned}\text{Cov}(\mathbf{x}, \mathbf{z}) &= E[(\mathbf{x} - \boldsymbol{\mu})\mathbf{z}^\top] \\ &= E[\mathbf{\Gamma}\mathbf{z}\mathbf{z}^\top] \\ &= \mathbf{\Gamma}\mathbf{A}\end{aligned}\tag{5.14}$$

oder

$$\text{Cov}(x_i, z_j) = \gamma_{ij}a_j \quad \text{für} \quad i, j = 1, \dots, p .\tag{5.15}$$

Somit ist die Korrelation zwischen Variablen und HK

$$\text{Corr}(x_i, z_j) = \lambda_{ij} = \frac{\text{Cov}(x_i, z_j)}{\sqrt{\text{Var}(x_i)}\sqrt{\text{Var}(z_j)}} = \frac{\gamma_{ij}a_j}{\sigma_{ii}^{\frac{1}{2}}a_j^{\frac{1}{2}}} = \gamma_{ij}\sqrt{\frac{a_j}{\sigma_{ii}}}\tag{5.16}$$

oder

$$\text{Corr}(\mathbf{x}, \mathbf{z}) = \mathbf{\Lambda} = \left(\text{Diag}(\mathbf{\Sigma})\right)^{-\frac{1}{2}}\mathbf{\Gamma}\mathbf{A}^{\frac{1}{2}} .\tag{5.17}$$

Bemerkung: Nach obigen Aussagen erkennt man sofort, dass die HK -Transformation nicht skaleninvariant ist und daher die Resultate von den Einheiten abhängen. Oft standardisiert man daher vorher die Variablen auf Mittel 0 und Varianz 1,

$$y_i = \frac{x_i - E(x_i)}{\sqrt{\text{Var}(x_i)}} \quad \text{für} \quad i = 1, \dots, p,\tag{5.18}$$

um Skaleninvarianz zu erreichen.

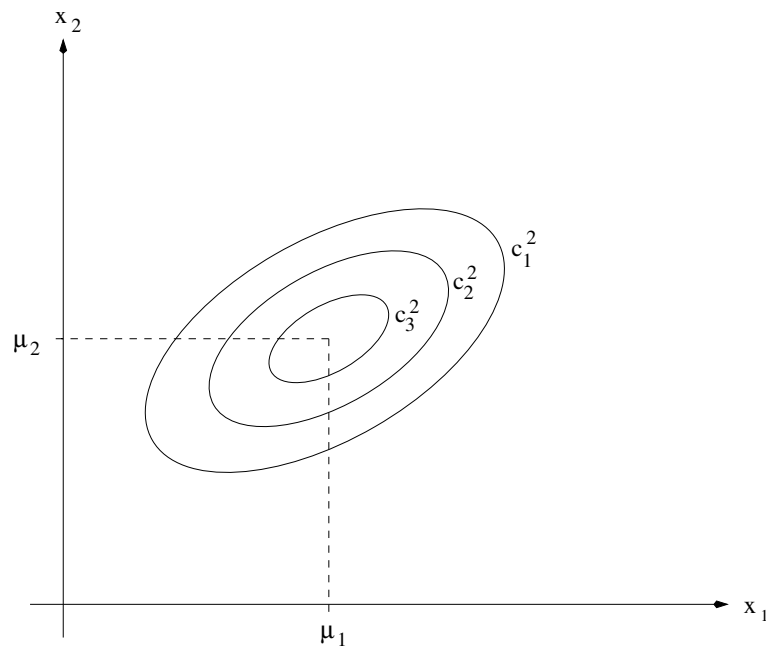
5.3 Geometrische Interpretation von Hauptkomponenten

In Kapitel 1.1 wurde die Mahalanobis-Distanz definiert. Setzt man, wie in (1.27) dieses Distanzmaß gleich einer Konstanten, also

$$(\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) := c^2 ,\tag{5.19}$$

so wird dadurch ein Ellipsoid im p -dimensionalen Raum mit Zentrum $\boldsymbol{\mu}$ beschrieben. Die Dichte der p -variaten Normalverteilung, die ja im Exponent die Mahalanobis-Distanz enthält, ist also auf solchen Ellipsoiden konstant.

Abbildung 5.1: Ellipsen mit konstanter Dichte bei der zweidimensionalen Normalverteilung



Durch die *HK*-Transformation

$$\mathbf{z} = \mathbf{\Gamma}^\top (\mathbf{x} - \boldsymbol{\mu}) \quad (5.20)$$

erhält man:

$$\begin{aligned} c^2 &= (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) = (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{\Gamma} \mathbf{A}^{-1} \mathbf{\Gamma}^\top (\mathbf{x} - \boldsymbol{\mu}) \\ &= \left[\mathbf{\Gamma}^\top (\mathbf{x} - \boldsymbol{\mu}) \right]^\top \mathbf{A}^{-1} \mathbf{\Gamma}^\top (\mathbf{x} - \boldsymbol{\mu}) = \mathbf{z}^\top \mathbf{A}^{-1} \mathbf{z} \\ &= \sum_{i=1}^p \frac{z_i^2}{a_i} = \frac{z_1^2}{a_1} + \frac{z_2^2}{a_2} + \dots + \frac{z_p^2}{a_p} \quad . \end{aligned}$$

Die Komponenten (z_1, \dots, z_p) von \mathbf{z} (d.h. die *HK*) repräsentieren also die Hauptachsen des Ellipsoids.

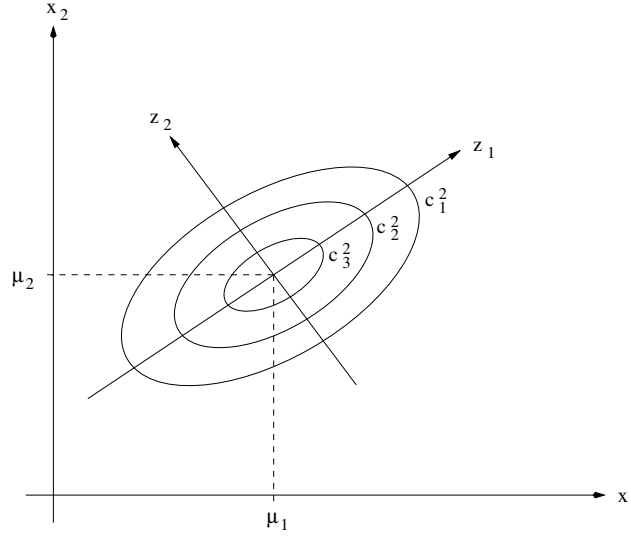
Die erste *HK* z_1 liegt in der größten Ausdehnung des Ellipsoids, die zweite *HK* z_2 liegt in der größten Ausdehnung, die orthogonal zu z_1 ist, die 3. *HK* z_3 liegt in der größten Ausdehnung, die orthogonal zu z_1 und z_2 ist usw. (siehe Abbildung 5.2).

5.4 Hauptkomponenten von Stichproben

Sei nun \mathbf{X} eine $(n \times p)$ Datenmatrix. Mittelvektor und Stichproben-Kovarianzmatrix können geschätzt werden durch

$$\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad , \quad (5.21)$$

Abbildung 5.2: Hauptkomponenten im zweidimensionalen Fall



\mathbf{x}_i ist die i -te Zeile von \mathbf{X} , und

$$\hat{\Sigma} = \mathbf{S} = \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top. \quad (5.22)$$

In Analogie zu (5.11) werden die *HK* der Stichprobe berechnet durch

$$\mathbf{Z} = (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top)\hat{\Gamma} \quad (5.23)$$

oder

$$\mathbf{z}_j = (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top)\hat{\gamma}_j \quad \text{für} \quad j = 1, \dots, p. \quad (5.24)$$

$\mathbf{1}$ ist ein Spaltenvektor bestehend aus n Einsen, und $\hat{\Gamma} = (\hat{\gamma}_1, \dots, \hat{\gamma}_p)$ ist die Matrix der Eigenvektoren von \mathbf{S} . Ähnlich zu den theoretischen *HK* erhalten wir die Zerlegung

$$\hat{\Gamma}^\top \mathbf{S} \hat{\Gamma} = \hat{\mathbf{A}}, \quad (5.25)$$

wobei $\hat{\mathbf{A}} = \text{Diag}(\hat{a}_1, \dots, \hat{a}_p)$ eine Diagonalmatrix der Eigenwerte (angeordnet in absteigender Größe) zu den zugehörigen Eigenvektoren von \mathbf{S} ist.

Die Matrix \mathbf{Z} der *HK* hat natürlich die selbe Dimension wie die Matrix \mathbf{X} . Die einzelnen *HK* \mathbf{z}_j kann man sich als Beobachtungen in einem orthogonal rotierten Raum vorstellen. Die Werte z_{ij} werden meist mit *Werte der HK* oder *scores* bezeichnet.

Analog zu (5.16) werden die Korrelationen zwischen Variablen und *HK* berechnet als

$$\hat{\lambda}_{ij} = \frac{\hat{\gamma}_{ij}\hat{a}_j^{\frac{1}{2}}}{s_{ii}^{\frac{1}{2}}} \quad \text{für} \quad i, j = 1, \dots, p, \quad (5.26)$$

wobei s_{ii} das i -te Diagonalelement von \mathbf{S} ist. In Matrixschreibweise wird die Korrelation durch

$$\hat{\mathbf{\Lambda}} = \left(\text{Diag}(\mathbf{S}) \right)^{-\frac{1}{2}} \hat{\mathbf{\Gamma}} \hat{\mathbf{A}}^{\frac{1}{2}} \quad (5.27)$$

berechnet.

Beispiel 5.4.1 Gegeben seien wieder die Prüfungsdaten aus Tabelle 2.1. Für diese Daten sollen nun die HK bestimmt werden.

Der Vektor der Spaltenmittel beträgt für diese Daten

$$\bar{\mathbf{x}} = (39.0 \quad 50.6 \quad 50.6 \quad 46.7 \quad 42.3)^\top .$$

Die empirische Kovarianzmatrix ergibt sich als

$$\mathbf{S} = \begin{pmatrix} 305.8 & 127.2 & 101.6 & 106.3 & 117.4 \\ & 172.8 & 85.2 & 94.7 & 99.0 \\ & & 112.9 & 112.1 & 121.9 \\ & & & 220.4 & 155.5 \\ & & & & 297.8 \end{pmatrix} .$$

Durch eine Spektralzerlegung von \mathbf{S} erhält man die Matrix der Eigenvektoren,

$$\hat{\mathbf{\Gamma}} = \begin{pmatrix} 0.51 & -0.75 & -0.30 & 0.30 & 0.08 \\ 0.37 & -0.21 & 0.42 & -0.78 & 0.19 \\ 0.35 & 0.08 & 0.15 & 0.00 & -0.92 \\ 0.45 & 0.30 & 0.60 & 0.52 & 0.29 \\ 0.53 & 0.55 & -0.60 & -0.18 & 0.15 \end{pmatrix} ,$$

sowie die Diagonalmatrix der Eigenwerte,

$$\hat{\mathbf{A}} = \text{Diag}(687.0, 202.1, 103.7, 84.6, 32.2).$$

Die HK werden dann berechnet durch

$$\mathbf{z}_j = (\mathbf{X} - \mathbf{1}\bar{\mathbf{x}}^\top) \hat{\mathbf{\gamma}}_j = (\mathbf{x}_1 \hat{\gamma}_{1j} + \dots + \mathbf{x}_p \hat{\gamma}_{pj}) - \mathbf{1} \bar{\mathbf{x}}^\top \hat{\mathbf{\gamma}}_j,$$

bzw. eingesetzt

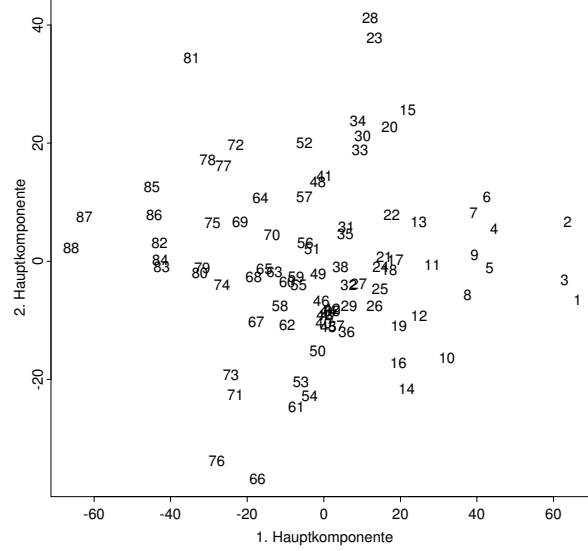
$$\begin{aligned} \mathbf{z}_1 &= 0.51\mathbf{x}_1 + 0.37\mathbf{x}_2 + 0.35\mathbf{x}_3 + 0.45\mathbf{x}_4 + 0.53\mathbf{x}_5 + 99.5 \cdot \mathbf{1} \\ \mathbf{z}_2 &= -0.75\mathbf{x}_1 - 0.21\mathbf{x}_2 + 0.08\mathbf{x}_3 + 0.30\mathbf{x}_4 + 0.55\mathbf{x}_5 + 1.4 \cdot \mathbf{1} \\ \mathbf{z}_3 &= -0.30\mathbf{x}_1 + 0.42\mathbf{x}_2 + 0.15\mathbf{x}_3 + 0.60\mathbf{x}_4 - 0.60\mathbf{x}_5 + 19.2 \cdot \mathbf{1} \\ \mathbf{z}_4 &= 0.30\mathbf{x}_1 - 0.78\mathbf{x}_2 + 0.00\mathbf{x}_3 + 0.52\mathbf{x}_4 - 0.18\mathbf{x}_5 - 11.5 \cdot \mathbf{1} \\ \mathbf{z}_5 &= 0.08\mathbf{x}_1 + 0.19\mathbf{x}_2 - 0.92\mathbf{x}_3 + 0.29\mathbf{x}_4 + 0.15\mathbf{x}_5 - 14.4 \cdot \mathbf{1} \end{aligned}$$

Die Varianzen der HK entsprechen den Eigenwerten,

$$\begin{aligned} \hat{a}_1 &= \widehat{\text{Var}}(\mathbf{z}_1) = 687.0 \\ \hat{a}_2 &= \widehat{\text{Var}}(\mathbf{z}_2) = 202.1 \\ \hat{a}_3 &= \widehat{\text{Var}}(\mathbf{z}_3) = 103.7 \\ \hat{a}_4 &= \widehat{\text{Var}}(\mathbf{z}_4) = 84.6 \\ \hat{a}_5 &= \widehat{\text{Var}}(\mathbf{z}_5) = 32.2 \quad . \end{aligned}$$

In Abbildung 5.3 sind die ersten beiden HK einander gegenübergestellt.

Abbildung 5.3: 1. und 2. Hauptkomponente der Prüfungsdaten



5.5 Anzahl der relevanten Hauptkomponenten

Eine Methode, die Anzahl k der relevanten *HK* festzustellen, ist ein Test auf Gleichheit der letzten (d.h. kleinsten) $(p - k)$ Eigenwerte, also $a_p = a_{p-1} = \dots = a_{k+1}$. Die letzten $(p - k)$ *HK* beschreiben somit gleiche Variation und enthalten daher gleich viel Information.

In der Praxis macht man oft eine Folge dieses Tests: Man startet mit $k = 0$ und erhöht k so lange, bis die Nullhypothese nicht mehr abgelehnt werden kann. Man kommt auf folgende Teststatistik:

$$\left(n - \frac{2p + 11}{6}\right) (p - k) \ln \left(\frac{m_a}{m_g}\right) \sim \chi^2_{(p-k+2)(p-k-1)/2} \quad (5.28)$$

mit

$$m_a = \frac{\hat{a}_{k+1} + \dots + \hat{a}_p}{p - k}, \quad m_g = \sqrt[p-k]{\hat{a}_{k+1} \dots \hat{a}_p}. \quad (5.29)$$

Wenn man mit standardisierten Variablen arbeitet, müssen die Eigenwerte von der Korrelationsmatrix $\boldsymbol{\rho}$ bestimmt werden. Eine Schätzung für $\boldsymbol{\rho}$ ist die Stichproben-Korrelationsmatrix \mathbf{R} . Ein Test für die Hypothese, daß die kleinsten $(p - k)$ Eigenwerte von \mathbf{R} gleich sind, ist dem obigen Test ähnlich. Die Teststatistik für diesen Fall ist

$$(n - 1)(p - k) \ln \left(\frac{m_a}{m_g}\right) \sim \chi^2_{(p-k+2)(p-k-1)/2}, \quad (5.30)$$

wobei m_a und m_g arithmetisches und geometrisches Mittel der letzten $(p - k)$ Eigenwerte von \mathbf{R} sind.

Beispiel 5.5.1 Gegeben seien wiederum die Prüfungsdaten. Zuerst untersuchen wir, ob die letzten vier Eigenwerte gleich sind, also

$$H_0 : \quad a_2 = a_3 = a_4 = a_5 .$$

Durch Anwendung der Formel für die Teststatistik erhält man eine χ^2_9 -Verteilung. Der Wert der Teststatistik ist 26.11 . Diese Nullhypothese ist signifikant zum 1%-Niveau und wird verworfen. Die neue Nullhypothese lautet nun

$$H_0 : \quad a_3 = a_4 = a_5 .$$

Man erhält eine χ^2_5 -Verteilung und der Wert der Teststatistik ist 7.30 . H_0 ist nicht einmal signifikant zum 5%-Niveau und wird daher angenommen.

5.5.1 „Faustregeln“

In der Literatur gibt es viele Regeln für die richtige Wahl der Anzahl von *HK*, die nicht statistisch begründet sind. Trotz dieses Mangels sind die Kriterien dafür weit verbreitet und werden auch in Statistik-Programmpaketen verwendet.

Eine Kenngröße, die die Entscheidung über die relevante Anzahl von *HK* erleichtern soll, ist der oben beschriebene Anteil der ersten k *HK* an der Gesamtvariation. Demnach sollte man so viele *HK* berücksichtigen, dass dieser Anteil zumindest 90% ist,

$$\frac{\sum_{j=1}^k \hat{a}_j}{\sum_{i=1}^p \hat{a}_i} \geq 90\% . \quad (5.31)$$

Ein anderes Kriterium besagt, dass man jene *HK* ausschließen soll, deren Eigenwerte geringer als der Durchschnitt sind.

Beispiel 5.5.2 Betrachten wir wieder unsere Prüfungsdaten, so berechnet man den Anteil jeder *HK* bzw. den Anteil der ersten j *HK* als

$\frac{\hat{a}_j}{\sum_i \hat{a}_i}$ in %	61.9	18.2	9.3	7.6	2.9
gesamt %	61.9	80.2	89.5	97.1	100

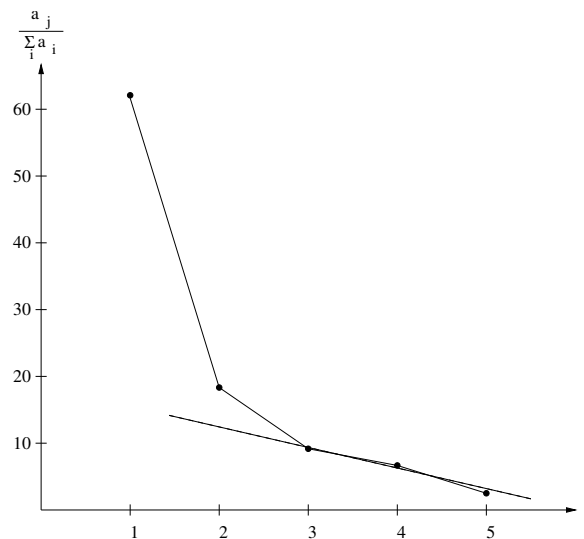
Nach dem 90%-Kriterium sollte man somit die ersten vier *HK* nehmen und nur die 5. Komponente weglassen. Das Kriterium mit dem Durchschnitt würde nur die erste Komponente einschließen.

Eine weitere Möglichkeit zum Bestimmen der maßgeblichen Variation bietet der *Scree-Graph* (‘‘Geröll’’-Grafik), der eine Gegenüberstellung vom Anteil jeder *HK* an der Gesamtvariation zur jeweiligen Nummer der *HK* darstellt. Jene Punkte, die auf einer Geraden liegen, werden weggelassen.

Beispiel 5.5.3 Im Fall der Prüfungsdaten ergibt sich der *Scree-Graph* in Abbildung 5.4.

Die Komponenten 3, 4 und 5 werden weggelassen.

Abbildung 5.4: Scree-Graph zu den Prüfungsdaten mit 5 HK



Eine andere Möglichkeit, den Anteil der Variation jeder Variable durch die *HK* zu erkennen, ist der durch (5.26) beschriebene Korrelationskoeffizient $\hat{\lambda}_{jk}$ der Variable \mathbf{x}_j mit der Komponente \mathbf{z}_k , bzw. das Quadrat dieser Zahl (Bestimmtheitsmaß). Sind die quadrierten Korrelationen mit den letzten *HK* nur mehr sehr gering, so können diese Komponenten ohne viel Informationsverlust weggelassen werden.

Beispiel 5.5.4 Bei den Prüfungsdaten ist der Anteil der letzten beiden *HK* an der Gesamtvariation 10.5%. Gefragt ist nun, wieviel Information verloren geht, wenn man diese Komponenten weglässt. Es ergeben sich folgende Werte für $\hat{\lambda}_{jk}^2$ (j sind die Variablen, k die Komponenten):

$\hat{\lambda}_{jk}^2$	\mathbf{z}_1	\mathbf{z}_2	\mathbf{z}_3	\mathbf{z}_4	\mathbf{z}_5
\mathbf{x}_1	0.574	0.371	0.030	0.024	0.001
\mathbf{x}_2	0.539	0.050	0.104	0.300	0.007
\mathbf{x}_3	0.727	0.018	0.019	0.000	0.243
\mathbf{x}_4	0.634	0.083	0.168	0.103	0.012
\mathbf{x}_5	0.660	0.204	0.125	0.009	0.002

Die erste *HK* erklärt 57.4% von \mathbf{x}_1 , 53.9% von \mathbf{x}_2 , 72.7% von \mathbf{x}_3 , 63.4% von \mathbf{x}_4 und 66.0% von \mathbf{x}_5 . Die ersten beiden *HK* erklären 94.5% von \mathbf{x}_1 etc. Man beachte, dass die letzte Komponente 24.3% der Variation von \mathbf{x}_3 , aber nur sehr wenig von den restlichen Variablen erklärt. Wenn man daher die letzte *HK* weglässt, verliert man viel mehr Information von der 3. Variablen als von den anderen. Wenn man die letzten beiden *HK* weglässt, verliert man die meiste Information von \mathbf{x}_2 und \mathbf{x}_3 , die wenigste von Variable 5.

In der Praxis wird oft ein Kompromiss zwischen verschiedenen Kriterien eingegangen.

Es gibt auch „modernere“ Verfahren zur Bestimmung der Anzahl von relevanten *HK* wie sogenannte *Resampling*-Methoden. Das sind nichtparametrische Schätzmethoden, die die Voraussetzung der Normalverteilung nicht benötigen (außer beim Testen). Die bekanntesten Resampling-Methoden sind *Jackknife* (Efron, 1982), *Bootstrap* (Efron, 1979; 1981) und *Cross Validation* (Stone, 1974; Eastment und Krzanowski, 1982).

5.6 Singulärwertzerlegung

Die Singulärwertzerlegung oder *singular value decomposition* (*SVD*) ist ein alternativer Algorithmus zur Berechnung der Hauptkomponenten, der keine Schätzung der Kovarianzmatrix benötigt. Dies hat vor allem Vorteile bei Daten mit $n < p$, wo die geschätzte Kovarianzmatrix immer singulär ist, und damit die Hauptkomponenten der Ordnung $> n$ automatisch Varianz 0 haben.

Wir gehen im folgenden davon aus, dass die Datenmatrix \mathbf{X} eine *zentrierte* $(n \times p)$ -Matrix mit reellen Zahlen ist. Es sind also die Spaltenmittel 0. Dann existiert eine orthogonale $(n \times n)$ -Matrix \mathbf{U} und eine orthogonale $(p \times p)$ -Matrix \mathbf{V} , sodass

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^\top \quad (5.32)$$

gilt. \mathbf{D} ist eine $(n \times p)$ -Matrix mit $d_{ii} \geq 0$ für $i = 1, \dots, \min(n, p)$, die restlichen Elemente sind null. Die positiven Werte d_{ii} werden *Singulärwerte von \mathbf{X}* genannt. Die Anzahl dieser positiven Werte entspricht dem Rang von \mathbf{X} .

Wenn der Rang von \mathbf{X} gleich $k < \min(n, p)$ ist, dann kann \mathbf{X} auch dargestellt werden durch

$$\mathbf{X} = \sum_{i=1}^k d_{ii} \mathbf{u}_i \mathbf{v}_i^\top. \quad (5.33)$$

\mathbf{u}_i bzw. \mathbf{v}_i bezeichnet die i -te Spalte von \mathbf{U} bzw. \mathbf{V} . Durch die Orthogonalität von \mathbf{U} und \mathbf{V} ergibt sich

$$\mathbf{X} \mathbf{X}^\top \mathbf{u}_i = d_{ii}^2 \mathbf{u}_i \quad (5.34)$$

und

$$\mathbf{X}^\top \mathbf{X} \mathbf{v}_i = d_{ii}^2 \mathbf{v}_i. \quad (5.35)$$

Daraus erkennt man, dass \mathbf{u}_i der i -te Eigenvektor von $\mathbf{X} \mathbf{X}^\top$ zum Eigenwert d_{ii}^2 , und \mathbf{v}_i der i -te Eigenvektor von $\mathbf{X}^\top \mathbf{X}$ zum selben Eigenwert d_{ii}^2 ist. Die Eigenwerte für $i = 1, \dots, k$ sind dabei größer als 0, die restlichen sind null.

Zusammenfassend heißt das, dass \mathbf{U} n orthogonale Eigenvektoren von $\mathbf{X} \mathbf{X}^\top$ in den Spalten hat, und dass \mathbf{V} p orthogonale Eigenvektoren von $\mathbf{X}^\top \mathbf{X}$ in den Spalten hat.

Man kann nun leicht den Zusammenhang zu einer Kovarianz-basierten Schätzung der Hauptkomponenten herstellen, so wie sie in Abschnitt 5.4 durchgeführt wurde. Wenn wir wiederum zentrierte Daten betrachten, dann waren die *HK* definiert als

$\mathbf{Z} = \mathbf{X}\hat{\mathbf{\Gamma}}$, und damit ist $\mathbf{X} = \mathbf{Z}\hat{\mathbf{\Gamma}}^\top$. Weiters ist die empirische Kovarianzmatrix in diesem Fall

$$\mathbf{S} = \frac{1}{n-1} \mathbf{X}^\top \mathbf{X}.$$

Die Matrix $\hat{\mathbf{\Gamma}}$ ist die Matrix der normierten Eigenvektoren von \mathbf{S} . Bei *SVD* ist \mathbf{V} die Matrix der normierten Eigenvektoren von $\mathbf{X}^\top \mathbf{X}$, woraus wir schließen können, dass $\hat{\mathbf{\Gamma}} \equiv \mathbf{V}$ ist. Aus Gleichung (5.35) schließen wir, dass $d_{ii}^2 = (n-1)\hat{a}_i$ ist. Wir haben somit die Darstellung

$$\mathbf{X} = \mathbf{Z}\mathbf{V}^\top = \mathbf{U}\mathbf{D}\mathbf{V}^\top, \quad (5.36)$$

und damit ist $\mathbf{Z} = \mathbf{U}\mathbf{D}$.

Mit diesen Überlegungen kann die Berechnung der Hauptkomponenten auch über eine andere Zielfunktion formuliert werden. Vorerst definieren wir noch die Frobenius Norm: Seien \mathbf{x}_i die Zeilen der Matrix \mathbf{X} , für $i = 1, \dots, n$. Dann ist die Frobenius Norm von \mathbf{X} definiert als

$$\|\mathbf{X}\|_F = \sqrt{\sum_{i=1}^n \|\mathbf{x}_i\|_2^2} = \sqrt{\sum_{i=1}^n \sum_{j=1}^p x_{ij}^2}$$

Laut oben haben wir

$$\mathbf{X}\mathbf{V} = \mathbf{U}\mathbf{D} = \mathbf{Z}.$$

Sei $\mathbf{V}_m = (\mathbf{v}_1, \dots, \mathbf{v}_m)$, und m kleiner als der Rang von \mathbf{X} . Dann stellen $\mathbf{X}\mathbf{V}_m$ genau die ersten m HK dar, was gleichbedeutend ist mit der Projektion von \mathbf{X} auf einen m -dimensionalen Unterraum. Man kann zeigen, dass

$$\mathbf{V}_m = \arg \max_{\mathbf{B}} \|\mathbf{X}\mathbf{B}\|_F^2$$

für eine beliebige $p \times m$ Matrix \mathbf{B} mit $\text{rank}(\mathbf{B}) \leq m$ und $\mathbf{B}^\top \mathbf{B} = \mathbf{I}$.

Eine äquivalente Formulierung ist folgende: Wir haben

$$\mathbf{X} = \mathbf{X}\mathbf{V}\mathbf{V}^\top = \mathbf{X}\mathbf{V}_m\mathbf{V}_m^\top + \mathbf{E}.$$

Dann gilt

$$\mathbf{V}_m = \arg \min_{\mathbf{B}} \|\mathbf{X} - \mathbf{X}\mathbf{B}\mathbf{B}^\top\|_F^2$$

mit gleicher Definition von \mathbf{B} wie oben. Man bemerke, dass $\hat{\mathbf{X}} := \mathbf{X}\mathbf{B}\mathbf{B}^\top$ die gleiche Dimension wie \mathbf{X} hat. Man kann somit $\hat{\mathbf{X}}$ als Rang m Approximation von \mathbf{X} auffassen, die in obigem Sinne optimal ist.

5.7 Biplots

Ein *Biplot* ist eine zweidimensionale grafische Darstellung von Objekten und Variablen in einer einzigen Abbildung. Biplots wurden von Gabriel (1971) entwickelt.

Das “Bi” in der Bezeichnung Biplot bezieht sich nicht auf die Dimension 2, sondern darauf, dass beide Größen (Objekte und Variablen) gleichzeitig repräsentiert werden. Selbstverständlich sind auch dreidimensionale Biplots denkbar, wenn auch etwas unpraktisch (Gabriel, 1986).

Wir werden hier Biplots kennenlernen, die den Zweck haben, Loadings und Scores von der Hauptkomponentenanalyse grafisch in einem Plot darzustellen. Es gibt auch andere Methoden von Biplots wie z.B. Multidimensionale Skalierung, Korrespondenzanalyse, Kanonische Biplots, Nichtlineare Biplots etc. Ein gutes Buch über diese Gebiete ist Gower und Hand (1996).

Zweidimensionale Plots sind gut handhabbar und übersichtlich. Eine Projektion der Daten auf zwei Dimensionen sollte allerdings auch voraussetzen, dass die Daten annähernd Rang 2 haben, dass also mit zwei Dimensionen der Großteil der Variabilität dargestellt werden kann. Für eine Datenmatrix mit höherem Rang sollten die ersten beiden Hauptkomponenten die Daten möglichst gut repräsentieren.

Sei nun \mathbf{X} eine zentrierte $(n \times p)$ -Matrix mit Rang k . Ein Biplot zeigt eine Darstellung von \mathbf{X} durch zwei Gruppen von Vektoren der Dimension n und p , die eine Approximation von \mathbf{X} mit Rang 2 sind. Diese Rang-2-Approximation bezeichnen wir mit $\mathbf{X}_{(2)}$. Eine Approximation von \mathbf{X} durch $\mathbf{X}_{(2)}$ im Sinne der Kleinsten Quadrate ist gegeben durch die ersten beiden Hauptkomponenten von \mathbf{X} , die entweder über die Schätzung der Kovarianzmatrix oder über SVD gewonnen werden können. Wir werden hier formal auf die Schätzung mit Singulärwertzerlegung eingehen. Wir benötigen also nur die ersten beiden Singulärwerte und die ersten beiden Spalten von \mathbf{U} und \mathbf{V} (nachfolgend trotzdem mit \mathbf{U} und \mathbf{V} bezeichnet):

$$\mathbf{X} \approx \mathbf{X}_{(2)} = \mathbf{U} \mathbf{D} \mathbf{V}^\top = (\mathbf{u}_1 \ \mathbf{u}_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix} \begin{pmatrix} \mathbf{v}_1^\top \\ \mathbf{v}_2^\top \end{pmatrix}. \quad (5.37)$$

$\mathbf{X}_{(2)}$ könnte aber auch anders aufgespalten werden, nämlich durch

$$\mathbf{X}_{(2)} = \mathbf{G} \mathbf{H}^\top \quad (5.38)$$

mit

$$\mathbf{G} = (\mathbf{u}_1 \ \mathbf{u}_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix}^{1-c} \quad (5.39)$$

und

$$\mathbf{H} = (\mathbf{v}_1 \ \mathbf{v}_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix}^c \quad (5.40)$$

für $0 \leq c \leq 1$. Je nach Wahl von c werden also die ersten beiden Singulärwerte auf die Matrizen \mathbf{G} und \mathbf{H} “aufgeteilt”. Ein Biplot besteht nun genau aus den Zeilen der Matrizen \mathbf{G} und \mathbf{H} , also aus $n + p$ zweidimensionalen Vektoren.

Für eine Wahl von $c = 0.5$ werden die Vektoren für Objekte und Variablen auf gleiche Skalierung gebracht.

Für $c = 1$ und Umskalierung,

$$\mathbf{G} = \begin{pmatrix} \mathbf{g}_{1.}^\top \\ \vdots \\ \mathbf{g}_{n.}^\top \end{pmatrix} = \sqrt{n-1} (\mathbf{u}_1 \ \mathbf{u}_2), \quad (5.41)$$

$$\mathbf{H} = \begin{pmatrix} \mathbf{h}_{1.}^\top \\ \vdots \\ \mathbf{h}_{p.}^\top \end{pmatrix} = \frac{1}{\sqrt{n-1}} (\mathbf{v}_1 \ \mathbf{v}_2) \begin{pmatrix} d_{11} & 0 \\ 0 & d_{22} \end{pmatrix}, \quad (5.42)$$

erhalten wir folgende Eigenschaften (Um die Notation zu vereinfachen, betrachten wir die Variablen als zentriert, also $\frac{1}{n} \sum_{i=1}^n x_{ij} = 0$ für $j = 1, \dots, p$):

- Das innere Produkt zwischen den Zeilen von \mathbf{G} und den Zeilen von \mathbf{H} approximieren die Datenmatrix \mathbf{X} ,

$$\mathbf{g}_{i.}^\top \mathbf{h}_{j.} = \sqrt{n-1} \mathbf{u}_{i.}^\top \frac{1}{\sqrt{n-1}} (\mathbf{v}_j^\top \mathbf{D})^\top = \mathbf{u}_{i.}^\top \mathbf{D} \mathbf{v}_j \approx x_{ij}. \quad (5.43)$$

- Die inneren Produkte zwischen den Zeilen von \mathbf{H} approximiert die Kovarianz,

$$\begin{aligned} \mathbf{H} \mathbf{H}^\top &= \left(\frac{1}{\sqrt{n-1}} \mathbf{V} \mathbf{D} \right) \left(\frac{1}{\sqrt{n-1}} \mathbf{D} \mathbf{V}^\top \right) = \frac{1}{n-1} \mathbf{V} \mathbf{D}^2 \mathbf{V}^\top \\ &= \frac{1}{n-1} (\mathbf{V} \mathbf{D} \mathbf{U}^\top) (\mathbf{U} \mathbf{D} \mathbf{V}^\top) = \frac{1}{n-1} \mathbf{X}_{(2)}^\top \mathbf{X}_{(2)} \\ &\approx \frac{1}{n-1} \mathbf{X}^\top \mathbf{X} = \mathbf{S}. \end{aligned} \quad (5.44)$$

Daraus folgt, dass die quadrierte Euklidische Norm der Zeilen von \mathbf{H} , $\|\mathbf{h}_{j.}\|^2 = \mathbf{h}_{j.}^\top \mathbf{h}_{j.}$, die Varianz approximiert. Darüber hinaus approximiert der Kosinus zwischen $\mathbf{h}_{i.}$ und $\mathbf{h}_{j.}$ ($i, j = 1, \dots, p$) die Korrelation zwischen den Variablen,

$$\cos(\mathbf{h}_{i.}, \mathbf{h}_{j.}) = \frac{\mathbf{h}_{i.}^\top \mathbf{h}_{j.}}{\|\mathbf{h}_{i.}\| \|\mathbf{h}_{j.}\|} \approx r_{ij}. \quad (5.45)$$

Als Konsequenz von (5.44) ergibt sich

$$\frac{1}{n-1} \sum_{l=1}^n (x_{li} - x_{lj})^2 = \frac{1}{n-1} (\mathbf{x}_{.i} - \mathbf{x}_{.j})^\top (\mathbf{x}_{.i} - \mathbf{x}_{.j}) \approx \|\mathbf{h}_{i.} - \mathbf{h}_{j.}\|^2. \quad (5.46)$$

D.h., dass die quadrierte Euklidische Distanz zwischen den Zeilen von \mathbf{H} die mittlere quadrierte Differenz zwischen den Variablen approximiert.

- Die Euklidische Distanz zwischen den Zeilen von \mathbf{G} approximiert die Mahalanobisdistanz zwischen den Beobachtungen,

$$\begin{aligned} \|\mathbf{g}_{i.} - \mathbf{g}_{j.}\|^2 &= (\mathbf{g}_{i.} - \mathbf{g}_{j.})^\top (\mathbf{g}_{i.} - \mathbf{g}_{j.}) = (n-1) (\mathbf{u}_{i.} - \mathbf{u}_{j.})^\top (\mathbf{u}_{i.} - \mathbf{u}_{j.}) \\ &\approx (\mathbf{x}_{i.} - \mathbf{x}_{j.})^\top \mathbf{S}^{-1} (\mathbf{x}_{i.} - \mathbf{x}_{j.}), \end{aligned} \quad (5.47)$$

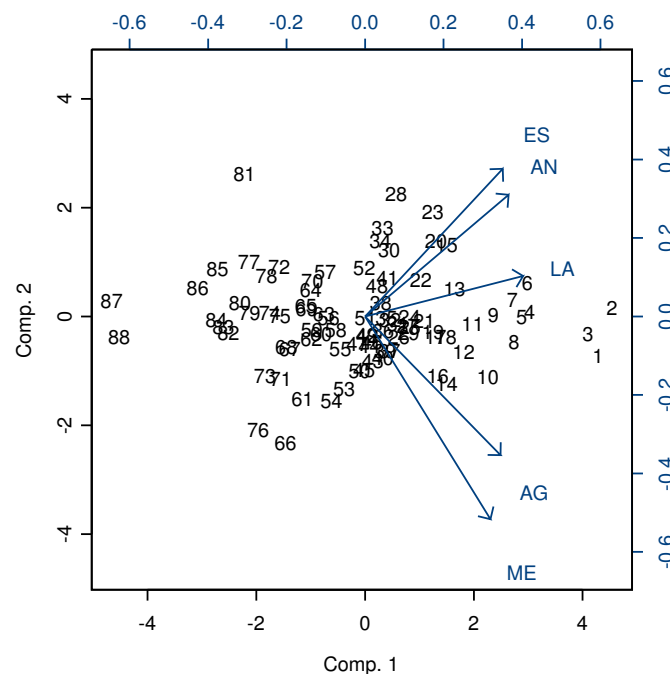
da

$$\mathbf{x}_{i.}^\top \mathbf{S}^{-1} \mathbf{x}_{j.} \approx (\mathbf{u}_{i.}^\top \mathbf{D} \mathbf{V}^\top) (n-1) (\mathbf{V} \mathbf{D}^{-2} \mathbf{V}^\top) (\mathbf{V} \mathbf{D} \mathbf{u}_{j.}) = (n-1) \mathbf{u}_{i.}^\top \mathbf{u}_{j.} \quad (5.48)$$

für $i, j = 1, \dots, n$.

Beispiel 5.7.1 Wir betrachten wiederum die Prüfungsdaten und erstellen damit einen Biplot (siehe Abbildung 5.5). Wir stellen fest, dass die Fächer ES und AN (und eventuell noch LA) in relativ starker Beziehung stehen, ebenso die Fächer AG und ME. Weiters ist ersichtlich, dass die x -Richtung zwischen guten und schlechten Resultaten generell unterscheidet, während die y -Richtung den Übergang von “mathematischen” zu “geometrischen” Fächern definiert. Die Einordnung jedes Studenten ist aus der Position im Plot (eigentlich aus der Projektion auf den Variablenvektor) leicht möglich. Z.B. ist Student 81 in den Fächern ME und AG sehr schlecht, bei den anderen Fächern liegt er im Mittelfeld.

Abbildung 5.5: Biplot der skalierten Prüfungsdaten



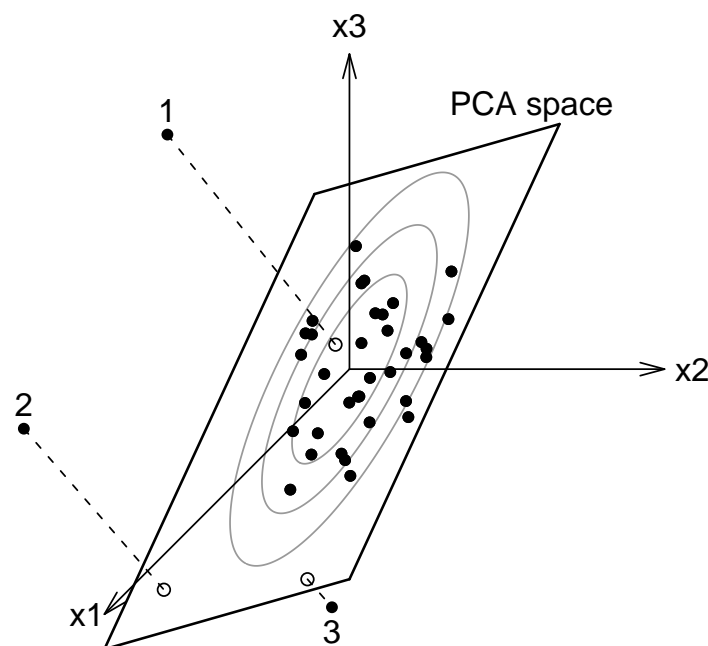
5.8 Diagnostik

Obwohl Hauptkomponentenanalyse keine typische Methodik zur multivariaten Ausreißererkennung darstellt, könnte man dennoch daran interessiert sein, ob einzelne Beobachtungen ungewöhnlich sind im Sinne von Abweichung von der Datenstruktur der Mehrheit der Daten. Dies wird durch zwei Distanz-Maße ermittelt, nämlich die *Score-Distanz* (SD) und die *orthogonale Distanz* (OD).

Abbildung 5.6 versucht dies, schematisch zu erläutern. Hier hat man Beobachtungen im 3-dimensionalen gegeben, und es wird der Raum der ersten beiden Hauptkomponenten visualisiert. Die (quadrierte) SD gibt die Mahalanobis-Distanz der Beobachtungen projiziert in diesen Raum an. Diese Distanz wird auch durch die

Ellipsen dargestellt. Die OD gibt die Distanz orthogonal zu diesem Raum an. Die meisten Beobachtungen sind im Raum der ersten beiden Hauptkomponenten gut dargestellt. Beobachtung 1 hat große OD, aber kleine SD. Beobachtung 2 hat große OD und große SD, und Beobachtung 3 hat große SD aber kleine OD. Analog zur Diagnostik in der Regression kann man hier Parallelen ziehen zu vertikalen Ausreißern (1), sowie guten (3) und schlechten (2) Hebelpunkten. Insbesondere schlechte Hebelpunkte können die klassische Schätzung der Hauptkomponenten stark beeinflussen. Daher ist diese Diagnostik erst sinnvoll für robust geschätzte Hauptkomponenten.

Abbildung 5.6: Diagnostik-Plot für Hauptkomponenten



Formal werden diese Distanzen folgendermaßen definiert. Sei $\hat{\mathbf{\Gamma}}_k = (\hat{\gamma}_1, \dots, \hat{\gamma}_k)$, also die Matrix der ersten k geschätzten *HK loadings*. Weiters sei $\mathbf{z}_i = (z_{i1}, \dots, z_{ik})^\top$ der i -te *score* Vektor der geschätzten *HK*, und $\hat{a}_1, \dots, \hat{a}_k$ die entsprechenden geschätzten *HK* Varianzen. Die Anzahl k von *HK* wird üblicherweise so gewählt, dass in etwa 80% der Varianz erklärt sind.

Die SD für die i -te Beobachtung ($i = 1, \dots, n$) ist definiert als

$$SD_i = \left(\sum_{j=1}^k \frac{z_{ij}^2}{\hat{a}_j} \right)^{1/2}.$$

Dies ist gleich mit der Mahalanobis-Distanz des *score* Vektors zum Zentrum der *HK* bezogen auf deren Kovarianzmatrix.

Die OD für die i -te Beobachtung \mathbf{x}_i ($i = 1, \dots, n$) ist definiert als

$$OD_i = \|\mathbf{x}_i - \hat{\Gamma}_k \mathbf{z}_i\|_2,$$

mit der euklidischen Norm $\|\cdot\|_2$, also als euklidische Distanz der Beobachtung zur Projektion in den Raum der ersten k HK.

Ähnlich wie bei multivariater Ausreißererkennung kann man Grenzen (cutoff-Werte) für beide Distanzen bestimmen, die auf ungewöhnlich hohe SD und OD hinweisen würden. Nachdem SD eine Mahalanobis-Distanz ist, ist ein geeigneter cutoff-Wert $\sqrt{\chi_{k;0.975}^2}$. Der cutoff-Wert für die OD ergibt sich über eine Approximation als

$$\left(\text{median}_i(OD_i^{2/3}) + \text{MAD}_i(OD_i^{2/3}) z_{0.975} \right)^{3/2},$$

wobei $z_{0.975}$ das 0.975-Quantil der $N(0, 1)$ ist. MAD steht für Median Absolute Deviation, und ist für univariate Werte y_1, \dots, y_n definiert als

$$\text{MAD} = 1.483 \cdot \text{median}_i(|y_i - \text{median}_j(y_j)|).$$

Beispiel 5.8.1 Für unsere Prüfungsdaten können diese Diagnostik-Plots sehr einfach konstruiert werden. Wir verwenden nachfolgend eine robuste Methode zur Schätzung der HK, implementiert als Funktion `PcaHubert()` im Paket `rrcov` – für Details siehe `help`-Seiten, sowie klassische HK mit `princomp()`. Der R Code ist folgend:

```
library(bootstrap)
data(scor) # Daten
library(rrcov)
res1 <- PcaHubert(scor, scale=TRUE) # robust
plot(res1)

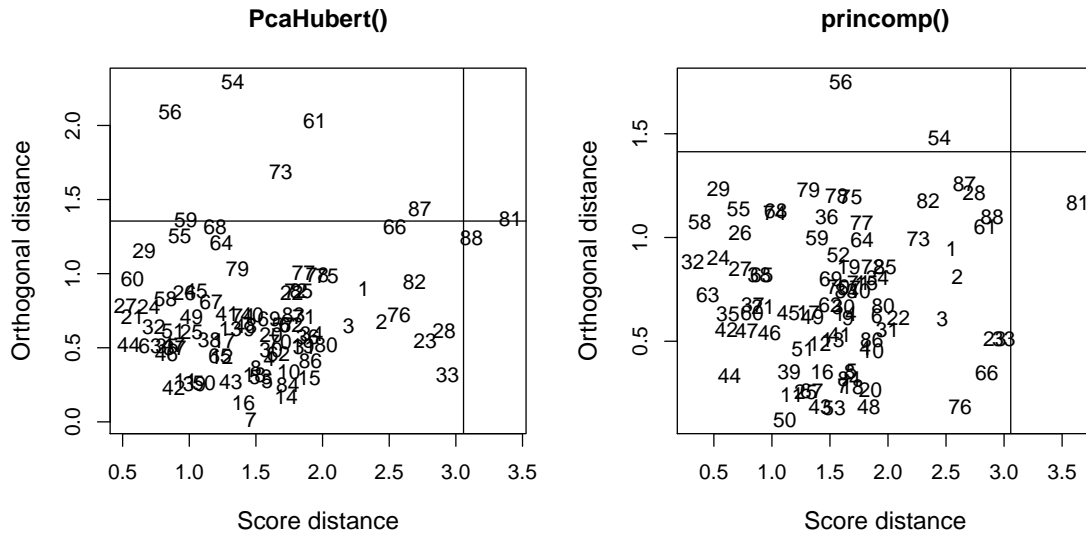
X <- scale(scor)
X.pca <- princomp(X) # nicht robust
library(chemometrics)
res <- pcaDiagplot(X, X.pca, a=3, plot=TRUE) # 3 HK
```

Die resultierenden Diagnostik-Plots sind in Abbildung 5.7 dargestellt. Die robuste Analyse (links) zeigt ein paar mehr ungewöhnliche Beobachtungen im Vergleich zur klassischen Analyse (rechts).

Literatur

- A. Basilevsky. *Statistical Factor Analysis and Related Methods: Theory and Applications*. Wiley & Sons, New York, 1994.
- T.F. Cox and A.A. Cox. *Multidimensional Scaling*. Chapman & Hall, London, 1994.
- H.T. Eastment and W.J. Krzanowski. Cross-validators choice of the number of components from a principal component analysis. *Technometrics*, 24:73–77, 1982.

Abbildung 5.7: Diagnostik-Plots für die Prüfungsdaten: links robust, rechts nicht robust



- B. Efron. Bootstrap Methods: Another Look at the Jackknife. *Ann. Statist.*, 7:1–26, 1979.
- B. Efron. Nonparametric Estimates of Standard Error: The Jackknife, the Bootstrap, and Other Methods. *Biometrika*, 68:589–599, 1981.
- B. Efron. *The Jackknife, the Bootstrap and Other Resampling Plans*. J.W. Arrowsmith Ltd., Bristol, 1982.
- K.R. Gabriel. The biplot graphic display of matrices with application to principal component analysis. *Biometrika*, 58(3):453–467, 1971.
- K.R. Gabriel and C.L. Odoroff. Use of three-dimensional biplots for diagnosis of models. In W. Gaul and M. Schader, editors, *Classification as a Tool of Research*, pages 153–159, 1986.
- J.C. Gower and D.J. Hand. *Biplots*. Chapman & Hall, London, 1996.
- J.C. Gower and S.A. Harding. Nonlinear biplots. *Biometrika*, 75(3):445–455, 1988.
- H. Hotelling. Analysis of a complex of statistical variables into principal components. *J. Educ. Psych.*, 24:417–441, 498–520, 1933.
- J.E. Jackson. *A User's Guide To Principal Components*. Wiley & Sons, New York, 1991.
- R. Johnson and D. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, London, 4th edition, 1998.
- K.V. Mardia, J.T. Kent, and J.M. Bibby. *Multivariate Analysis*. Acad. Press, London, 1979.
- K. Pearson. On lines and planes of closest fit to systems of points in space. *Phil. Mag. (6)*, 2:559–572, 1901.
- M. Stone. Cross-validatory choice and assessment of statistical prediction. *J. R. Stat. Soc. B*, 36:111–133, 1974.

Kapitel 6

Faktorenanalyse

6.1 Einleitung

Betrachtet man an n Objekten (Personen) p Merkmale (Körpergröße, Gewicht, Alter, ...), so werden die Merkmale in der Regel voneinander abhängen, d.h. auch miteinander korrelieren. Es ist anzunehmen, dass sie einander wechselseitig beeinflussen, oder dass eine dritte Größe, die man nicht direkt messen kann, den Einfluss bestimmt (z.B. Statur).

Das faktorenanalytische Modell geht immer davon aus, dass das Messbare nur eine Erscheinungsform von Größen ist, die im Hintergrund stehen, und die man direkt nicht messen kann. In vielen Fällen ist diese Annahme realistisch. Es ist daher von Interesse, ob sich aus den beobachteten Merkmalen eine Größe isolieren lässt, ein sogenannter Faktor, der die Zusammenhänge erklärt.

Das Hauptziel der Faktorenanalyse ist die Ableitung hypothetischer Größen oder Faktoren aus einer Menge beobachteter Merkmale. Die Faktoren sollen möglichst einfach sein und die Beobachtungen hinreichend genau (interpretierbar) beschreiben und erklären.

Der klassische Bereich der Faktorenanalyse ist die Theorie der Intelligenz. Man kann zunächst annehmen, dass sich Intelligenz aus zahlreichen Elementarfähigkeiten wie Gedächtnis, Aufmerksamkeit, Urteilsfähigkeit usw. zusammensetzt. Solche Größen versucht man durch psychologische Testverfahren zu messen, um dann mit Hilfe der Faktorenanalyse den Faktor „Intelligenz“ zu isolieren. Anhand dieses Faktors kann die beobachtete Personengruppe leichter unterteilt werden als auf Grund vieler einzelner Testergebnisse.

6.2 Das Faktorenmodell

6.2.1 Definition

Sei $\mathbf{x} = (x_1, x_2, \dots, x_p)^\top$ ein stochastischer Vektor, x_1, x_2, \dots, x_p die Merkmale bzw. Variablen.

Sei weiters $\mathbf{y} = (y_1, y_2, \dots, y_p)^\top$ der Vektor, der durch

$$y_i = \frac{x_i - E(x_i)}{\sqrt{Var(x_i)}}, \quad i = 1, \dots, p, \quad ,$$

auf Mittel Null und Varianz Eins standardisierten Merkmale.

Es wird angenommen, dass die Elemente von \mathbf{y} , abgesehen von einem Fehlerterm \mathbf{e} , durch eine kleinere Anzahl $k < p$ von unbekannten Zufallsvariablen $\mathbf{f} = (f_1, \dots, f_k)^\top$ dargestellt werden können. Dann lässt sich das k -Faktorenmodell folgendermaßen angeben:

$$\mathbf{y} = \mathbf{\Lambda} \mathbf{f} + \mathbf{e} \quad . \quad (6.1)$$

Dabei ist $\mathbf{\Lambda} = (\lambda_{ij})$ eine $(p \times k)$ -Matrix mit konstanten Werten. Sie wird *Ladungsmatrix* genannt und beschreibt den Zusammenhang zwischen Faktoren und Merkmalen.

Je nach Größe und Verteilung der *Ladungen* λ_{ij} ($i = 1, \dots, p$) spricht man bei f_j von einem

- *allgemeinen Faktor* (engl. *general factor*), wenn fast *alle* der oben angeführten Ladungen deutlich von Null verschieden sind;
- *gemeinsamen Faktor* (engl. *common factor*), wenn *mindestens zwei* Ladungen deutlich von Null verschieden sind.

Der Fehlerterm $\mathbf{e} = (e_1, \dots, e_p)^\top$ wird *Einzelrestkomponente* (engl. *unique factor*) genannt.

Bei diesem Modell wird folgendes vorausgesetzt:

$$\begin{aligned} E(\mathbf{f}) &= \mathbf{0}, \quad Cov(e_i, e_j) = 0 \quad (i \neq j), \\ E(\mathbf{e}) &= \mathbf{0}, \quad Cov(\mathbf{f}, \mathbf{e}) = \mathbf{0}, \\ Var(f_i) &= 1. \end{aligned}$$

Die Kovarianzmatrix von \mathbf{e} lässt sich somit in der Form

$$Cov(\mathbf{e}) = \mathbf{\Psi} = Diag(\psi_{11}, \dots, \psi_{pp}) \quad (6.2)$$

schreiben.

Für die Korrelationsmatrix $\boldsymbol{\rho} = [(\rho_{ij})]_{i,j=1,\dots,p}$ von \mathbf{x} gilt nun:

$$\begin{aligned} \boldsymbol{\rho} &= Corr(\mathbf{x}) = Cov(\mathbf{y}) = Cov(\mathbf{\Lambda} \mathbf{f} + \mathbf{e}) \\ &= \underbrace{\mathbf{\Lambda} Cov(\mathbf{f}) \mathbf{\Lambda}^\top}_{\mathbf{\Phi}} + \underbrace{\mathbf{\Lambda} Cov(\mathbf{f}, \mathbf{e})}_{\mathbf{0}} + \underbrace{Cov(\mathbf{e}, \mathbf{f}) \mathbf{\Lambda}^\top}_{\mathbf{0}} + \underbrace{Cov(\mathbf{e})}_{\mathbf{\Psi}} \\ &= \mathbf{\Lambda} \mathbf{\Phi} \mathbf{\Lambda}^\top + \mathbf{\Psi} \quad . \end{aligned}$$

$\mathbf{\Phi}$ ist die $(k \times k)$ -Korrelationsmatrix der Faktoren.

Fordert man zusätzlich, dass die Faktoren unkorreliert sind, also

$$Cov(\mathbf{f}) = \mathbf{\Phi} = \mathbf{I} \quad , \quad (6.3)$$

so erhält man

$$\boldsymbol{\rho} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi} \quad (6.4)$$

bzw.

$$\boldsymbol{\rho}_{red} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top = \begin{pmatrix} \kappa_1^2 & \rho_{12} & \cdots & \rho_{1p} \\ \rho_{21} & \kappa_2^2 & \cdots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \cdots & \kappa_p^2 \end{pmatrix} = \boldsymbol{\rho} - \boldsymbol{\Psi} \quad (6.5)$$

$\boldsymbol{\rho}_{red}$ bezeichnet die *reduzierte Korrelationsmatrix*. Die Diagonalelemente $\kappa_i^2 = 1 - \psi_{ii} = \sum_{j=1}^k \lambda_{ij}^2$ ($i = 1, \dots, p$) werden *Kommunalitäten* genannt. Sie entsprechen den Zeilensummen der quadrierten Faktorenladungen und beschreiben den Anteil der Varianz von y_i , der durch die gemeinsamen Faktoren erklärt wird.

6.2.2 Nichteindeutigkeit der Faktorenladungen

Sei \mathbf{G} eine orthogonale ($k \times k$)-Matrix. Dann gilt wegen $\mathbf{G}^{-1} = \mathbf{G}^\top$

$$\mathbf{y} = (\boldsymbol{\Lambda}\mathbf{G})(\mathbf{G}^\top \mathbf{f}) + \mathbf{e} \quad (6.6)$$

Nachdem die neuen Faktoren $\mathbf{G}^\top \mathbf{f}$ die Modellannahmen erfüllen,

$$\begin{aligned} E(\mathbf{G}^\top \mathbf{f}) &= \mathbf{0}, & Cov(\mathbf{G}^\top \mathbf{f}) &= \mathbf{I} \\ \text{und} & & Cov(\mathbf{G}^\top \mathbf{f}, \mathbf{e}) &= \mathbf{0} \end{aligned}$$

ist das k -Faktorenmodell auch gültig mit den neuen Faktoren, d.h.

$$\boldsymbol{\rho} = (\boldsymbol{\Lambda}\mathbf{G})(\mathbf{G}^\top \boldsymbol{\Lambda}^\top) + \boldsymbol{\Psi} \quad (6.7)$$

Das bedeutet aber, dass die Faktorenladungsmatrix nicht eindeutig bestimmt ist. Um eine eindeutige Lösung für die Faktorenladungsmatrix zu erhalten, fordern wir zusätzlich, dass entweder $\boldsymbol{\Lambda}^\top \boldsymbol{\Psi}^{-1} \boldsymbol{\Lambda}$ oder $\boldsymbol{\Lambda}^\top \boldsymbol{\Lambda}$ eine Diagonalmatrix sein soll.

6.2.3 Parameterschätzung

In der Praxis ist eine Datenmatrix aus beobachteten Werten gegeben. Daraus kann die Korrelationsmatrix $\hat{\boldsymbol{\rho}}$ geschätzt werden. Anhand von $\hat{\boldsymbol{\rho}}$ sollen nun Schätzungen $\hat{\boldsymbol{\Lambda}}$ und $\hat{\boldsymbol{\Psi}}$ für $\boldsymbol{\Lambda}$ und $\boldsymbol{\Psi}$ gefunden werden, die sowohl

$$\hat{\boldsymbol{\Lambda}}^\top \hat{\boldsymbol{\Psi}}^{-1} \hat{\boldsymbol{\Lambda}} = \text{Diag} \quad \text{oder} \quad \hat{\boldsymbol{\Lambda}}^\top \hat{\boldsymbol{\Lambda}} = \text{Diag} \quad (6.8)$$

als auch

$$\hat{\boldsymbol{\rho}} \cong \hat{\boldsymbol{\Lambda}} \hat{\boldsymbol{\Lambda}}^\top + \hat{\boldsymbol{\Psi}} \quad (6.9)$$

erfüllen.

Es stellt sich die Frage, ob das Faktorenmodell eine einfachere Interpretation der Daten liefert als die Korrelationsmatrix. Dazu kann die Anzahl der Parameter im uneingeschränkten Modell mit der Anzahl der freien Parameter im k -Faktorenmodell

verglichen werden. $\mathbf{\Lambda}$ und $\mathbf{\Psi}$ bestimmen gemeinsam $pk + p$ freie Parameter. Die Forderung, dass $\mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{\Lambda}$ oder $\mathbf{\Lambda}^\top \mathbf{\Lambda}$ eine Diagonalmatrix ist, führt zu $\frac{1}{2}k(k-1)$ Einschränkungen. Es bleiben daher $pk + p - \frac{1}{2}k(k-1)$ Parameter. Die Anzahl der freien Parameter im uneingeschränkten Modell beträgt $\frac{1}{2}p(p+1)$. Die Differenz s ist somit

$$\begin{aligned} s &= \frac{1}{2}p(p+1) - \left(pk + p - \frac{1}{2}k(k-1) \right) \\ &= \frac{1}{2}(p-k)^2 - \frac{1}{2}(p+k) \quad . \end{aligned}$$

Je nach dem Wert von s können drei Fälle unterschieden werden:

1. $s < 0$: Es gibt mehr Parameter als Gleichungen. Daher können unendlich viele Lösungen für $\hat{\mathbf{\Lambda}}$ und $\hat{\mathbf{\Psi}}$ gefunden werden. Das Faktorenmodell ist somit nicht wohldefiniert.
2. $s = 0$: Es gibt gleich viele Parameter wie Gleichungen. $\hat{\mathbf{\Lambda}}$ und $\hat{\mathbf{\Psi}}$ können zwar eindeutig bestimmt werden, jedoch stellt das Faktorenmodell keine einfachere Beschreibung der Daten dar.
3. $s > 0$: Es gibt mehr Gleichungen als Parameter. $\hat{\mathbf{\Lambda}}$ und $\hat{\mathbf{\Psi}}$ können nicht exakt bestimmt werden, jedoch liefert das Faktorenmodell eine einfachere Interpretation der Daten als die Korrelationsmatrix.

Im folgenden wird nur der Fall $s > 0$ betrachtet.

6.3 Vorgangsweise

Der wesentliche Schritt in der Faktorenanalyse besteht in der Schätzung der Ladungsmatrix $\hat{\mathbf{\Lambda}}$. Dazu gibt es verschiedene Verfahren, wie z.B. die Hauptkomponenten- bzw. Hauptfaktorenanalyse, die Maximum-Likelihood-Methode, die Zentroidmethode usw.

In diesem Abschnitt werden nur die Hauptfaktorenanalyse und die Maximum-Likelihood-Methode behandelt. Die Hauptkomponentenanalyse entspricht der Ermittlung von Hauptkomponenten der *nicht* reduzierten Korrelationsmatrix. Mit der Zentroidmethode werden die Achsen durch sogenannte *Schwerpunkte* gelegt.

6.3.1 Hauptfaktorenanalyse

Ausgehend von dem Modell

$$\boldsymbol{\rho} = \mathbf{\Lambda} \mathbf{\Lambda}^\top + \mathbf{\Psi} \tag{6.10}$$

werden zuerst die Kommunalitäten geschätzt. Anschließend schätzt man die Ladungsmatrix.

Kommunalitätenschätzung

Die Kommunalitäten

$$\kappa_i^2 = \sum_{j=1}^k \lambda_{ij}^2 = 1 - \psi_{ii} \quad (i = 1, \dots, p) \quad (6.11)$$

beschreiben den Anteil der Varianz, der durch das gemeinsame Faktorenmodell erklärt wird. Sie können nur Werte zwischen 0 und 1 annehmen. Es gibt mehrere Möglichkeiten, die Kommunalitäten zu schätzen:

1. Der höchste Korrelationskoeffizient : $\max_{i \neq j} |\rho_{ij}|$
Für die Schätzung der Kommunalitäten wird in jeder Spalte von $\boldsymbol{\rho}$ das betragsmäßig größte Nichtdiagonalelement herangezogen. Diese Methode ist nur bei größerer Merkmalsanzahl empfehlenswert, da das Ergebnis zufallsabhängig ist.
2. Das Quadrat des multiplen Korrelationskoeffizienten: $\rho_{i,12\dots i(\dots p)}^2$
Aus der Regressionsrechnung ist es als SMC (engl. *squared multiple correlation*) bekannt und beschreibt den Varianzanteil des i-ten Merkmals erklärt durch die anderen. Das Quadrat des multiplen Korrelationskoeffizienten des i-ten Merkmals kann folgendermaßen berechnet werden:

$$\rho_{i,12\dots i(\dots p)}^2 = 1 - \frac{1}{\rho^{ii}} \quad (6.12)$$

Dabei ist ρ^{ii} das i-te Diagonalelement von $\boldsymbol{\rho}^{-1}$.

3. Die Iteration
Bei der Iteration muss zuerst eine Entscheidung über die Anzahl k der zu extrahierenden Faktoren getroffen werden. Dazu kann man nach einem in der Hauptkomponentenanalyse besprochenen Kriterium vorgehen.

Vorerst werden Schätzungen der Kommunalitäten (SMC oder Werte zwischen 0 und 1) in die Diagonale von $\boldsymbol{\rho}$ eingesetzt. Aus dieser reduzierten Korrelationsmatrix werden k Faktoren extrahiert (siehe Ende dieses Abschnitts). Aus dem Faktorenmuster (Ladungsmatrix) werden die neuen Kommunalitäten

$$\kappa_{i_{neu}}^2 = \sum_{j=1}^k \lambda_{ij}^2 \quad \text{für } i = 1, \dots, p \quad (6.13)$$

berechnet und als neue Schätzungen in die Diagonale von $\boldsymbol{\rho}$ eingesetzt. Die Iteration wird abgebrochen, wenn sich die Kommunalitäten nur mehr geringfügig ändern.

Die Konvergenz dieses Verfahrens ist theoretisch nicht nachgewiesen. Außerdem wurde nicht gezeigt, unter welchen Bedingungen dieses Verfahren konvergiert und ob der Grenzwert mit den tatsächlichen Kommunalitäten übereinstimmt. Dennoch ist diese Methode am empfehlenswertesten.

Wählt man die Kommunalitäten zu hoch, so wird ein Teil der Einzelvarianz in die gemeinsame Varianz hineingepresst und das Faktorenmuster dadurch verändert. Wählt man sie zu niedrig, dann geht gemeinsame Varianz für die Bestimmung der Faktoren verloren.

Bei einer größeren Anzahl von Merkmalen ist der Einfluss ungenauer Kommunalitätsschätzungen auf die Lösung relativ gering. Bei kleinen Variablenanzahlen kann jedoch die Kommunalitätsschätzung entscheidend sein.

Schätzung der Ladungsmatrix

Nachdem nun die Kommunalitäten bekannt sind, kann man für die Schätzungen der Ladungen von der reduzierten Korrelationsmatrix

$$\boldsymbol{\rho}_{red} = \boldsymbol{\rho} - \boldsymbol{\Psi} \quad (6.14)$$

ausgehen, bei der in der Diagonale anstelle der Einsen die Kommunalitäten stehen.

Nach dem Spektralzerlegungssatz gilt

$$\boldsymbol{\rho} - \boldsymbol{\Psi} = \boldsymbol{\Gamma} \mathbf{A} \boldsymbol{\Gamma}^\top = \sum_{i=1}^p a_i \boldsymbol{\gamma}_i \boldsymbol{\gamma}_i^\top, \quad (6.15)$$

wobei $\boldsymbol{\Gamma} = (\boldsymbol{\gamma}_{.1}, \dots, \boldsymbol{\gamma}_{.p})$ sowie $\mathbf{A} = \text{Diag}(a_1, \dots, a_p)$ gilt und $\boldsymbol{\gamma}_i$ bzw. a_i die Eigenvektoren bzw. die Eigenwerte von $\boldsymbol{\rho} - \boldsymbol{\Psi}$ sind ($i = 1, \dots, p$).

Das Faktorenmodell ist somit

$$\boldsymbol{\rho} - \boldsymbol{\Psi} = \boldsymbol{\Lambda} \boldsymbol{\Lambda}^\top = \boldsymbol{\Gamma} \mathbf{A} \boldsymbol{\Gamma}^\top. \quad (6.16)$$

Für die Schätzung der Faktorenladungsmatrix gilt deshalb

$$\hat{\boldsymbol{\Lambda}} = \boldsymbol{\Gamma} \mathbf{A}^{1/2}. \quad (6.17)$$

Werden k Faktoren extrahiert ($1 \leq k \leq p$), so kann diese Gleichung geschrieben werden als

$$\hat{\boldsymbol{\lambda}}_{.j} = a_j^{1/2} \boldsymbol{\gamma}_{.j} \quad \text{für } j = 1, \dots, k. \quad (6.18)$$

Für die Schätzung der merkmalseigenen Varianzen gilt

$$\hat{\psi}_{ii} = 1 - \sum_{j=1}^k \hat{\lambda}_{ij}^2 \quad \text{für } i = 1, \dots, p. \quad (6.19)$$

Die Lösung ist zulässig, wenn alle $\hat{\psi}_{ii} \geq 0$ sind.

Beispiel 6.3.1 Gegeben seien die Prüfungsdaten aus Tabelle 2.1. Die empirische Korrelationsmatrix ergibt sich als

$$\mathbf{R} = \begin{pmatrix} 1 & 0.553 & 0.547 & 0.410 & 0.389 \\ & 1 & 0.610 & 0.485 & 0.437 \\ & & 1 & 0.711 & 0.665 \\ & & & 1 & 0.607 \\ & & & & 1 \end{pmatrix}.$$

Wählt man die Anzahl der Faktoren k größer als 2, so ist das Faktorenmodell nicht wohldefiniert ($k = 3 \Rightarrow s = -2 < 0$). Das bedeutet, dass entweder ein oder zwei Faktoren bestimmt werden können.

Zum Vergleich wurden die Kommunalitäten sowohl mit dem Quadrat des multiplen Korrelationskoeffizienten (SMC) als auch mit dem höchsten Korrelationskoeffizienten $\max_{i \neq j} |r_{ij}|$ geschätzt.

Für die Inverse der Korrelationsmatrix ergibt sich

$$\mathbf{R}^{-1} = \begin{pmatrix} 1.603 & -0.558 & -0.510 & 0.001 & -0.041 \\ & 1.802 & -0.0659 & -0.152 & -0.039 \\ & & 3.047 & -1.113 & -0.864 \\ & & & 2.178 & -0.515 \\ & & & & 1.921 \end{pmatrix}.$$

Aus diesem Ergebnis wird nun der SMC-Schätzer berechnet. Somit ergibt sich für die Schätzungen der Kommunalitäten:

Schätzer	$\hat{\kappa}_1^2$	$\hat{\kappa}_2^2$	$\hat{\kappa}_3^2$	$\hat{\kappa}_4^2$	$\hat{\kappa}_5^2$
SMC	0.376	0.445	0.672	0.541	0.479
$\max_{i \neq j} r_{ij} $	0.553	0.610	0.711	0.711	0.665

In diesem Beispiel werden die Kommunalitäten mit dem höchsten Korrelationskoeffizienten geschätzt. Aus der reduzierten Korrelationsmatrix

$$\mathbf{R} - \Psi = \begin{pmatrix} 0.553 & 0.553 & 0.547 & 0.410 & 0.389 \\ & 0.610 & 0.610 & 0.485 & 0.437 \\ & & 0.711 & 0.711 & 0.665 \\ & & & 0.711 & 0.607 \\ & & & & 0.665 \end{pmatrix}$$

werden die Eigenwerte berechnet:

$$a_1 = 2.84 \quad a_2 = 0.38 \quad a_3 = 0.08 \quad a_4 = 0.02 \quad a_5 = -0.05$$

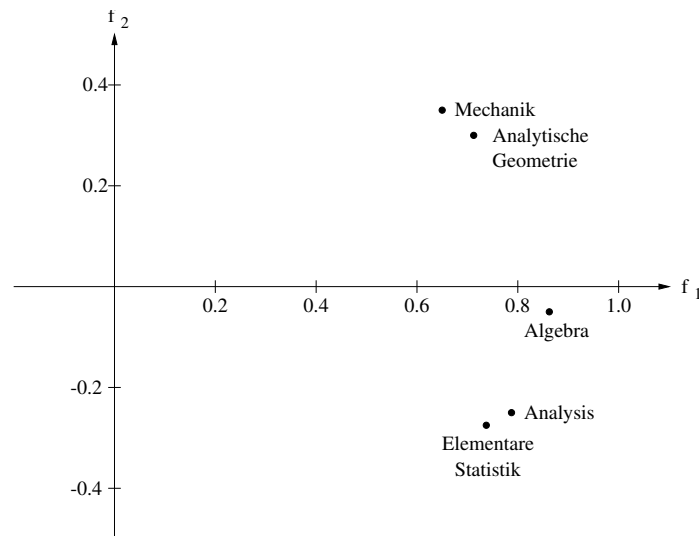
Aus der Größe der Eigenwerte kann man erkennen, dass die Annahme eines 1- bzw. 2-Faktorenmodells sinnvoll war. Für die Faktorenschätzungen und Kommunalitäten des 1- bzw. 2-Faktorenmodells ergibt sich:

Merkmale	$k=1$		$k=2$		
	κ_i^2	λ_1	κ_i^2	λ_1	λ_2
ME	0.417	0.646	0.543	0.646	0.354
AG	0.506	0.711	0.597	0.711	0.303
LA	0.746	0.864	0.749	0.864	-0.051
AN	0.618	0.786	0.680	0.786	-0.249
ES	0.551	0.742	0.627	0.742	-0.276

Beim 1-Faktorenmodell kann man erkennen, dass die Ladungen in etwa gleich groß sind. In der „Ladungsmatrix“ ist daher keine Struktur ersichtlich. Anders ist das

beim 2-Faktorenmodell. Während sich die erste Spalte der Ladungsmatrix nicht ändert, zeigt die zweite Spalte positive Ladung bei den Merkmalen 1 und 2, keine Ladung bei Merkmal 3 und negative Ladung bei den Merkmalen 4 und 5. Der erste Faktor repräsentiert somit die durchschnittliche Leistung in den Prüfungsfächern, der zweite weniger wichtige Faktor ($a_2 = 0.38 \ll 2.84 = a_1$) beschreibt den Unterschied zwischen den Fächern (siehe auch Abbildung 6.1).

Abbildung 6.1: Faktorenladungen



Es muss darauf hingewiesen werden, dass alle Kommunalitäten wesentlich kleiner als 1 sind, d.h. es bleibt ein bedeutender Anteil der Varianz jeder Variablen unerklärt.

6.3.2 Maximum-Likelihood-Methode

Diese Methode wurde von Lawley (1940, 1941) eingeführt. Es wird angenommen, dass die betrachteten Variablen normalverteilt sind und dass die gemeinsamen Faktoren orthogonal sind. Wie auch in der Hauptfaktorenanalyse wird das Modell

$$\boldsymbol{\rho} = \boldsymbol{\Lambda}\boldsymbol{\Lambda}^\top + \boldsymbol{\Psi} \quad (6.20)$$

vorausgesetzt. Die Idee dieses Verfahrens ist, durch Maximierung der logarithmierten Plausibilitätsfunktion konsistente und effiziente Schätzungen der unbekannten Parameter $\boldsymbol{\Lambda}$ und $\boldsymbol{\Psi}$ zu erhalten.

Setzt man die SP-Korrelationsmatrix \mathbf{R} als Realisierung für den Schätzer $\hat{\boldsymbol{\rho}}$ ein, so erhält man die logarithmierte Plausibilitätsfunktion

$$l = -\frac{1}{2}n \ln |\boldsymbol{\rho}| - \frac{1}{2}n \operatorname{tr} \boldsymbol{\rho}^{-1} \mathbf{R} \quad (6.21)$$

l soll maximiert werden, d.h. die partiellen Ableitungen

$$\frac{\partial l}{\partial \lambda_{ij}} = -n \left(\sum_t \lambda_{tj} \rho^{ti} - \sum_{t,u,v} \lambda_{tj} \rho^{tu} r_{uv} \rho^{vi} \right) \quad (6.22)$$

und

$$\frac{\partial l}{\partial \psi_{ii}} = -\frac{1}{2}n\left(\rho^{ii} - \sum_{t,u} \rho^{it} r_{tu} \rho^{ui}\right) \quad (6.23)$$

werden Null gesetzt ($i = 1, \dots, p$).

Die geschätzte Ladungsmatrix $\hat{\Lambda}$, die merkmalseigenen Varianzen $\hat{\Psi}$ und damit die geschätzte Korrelationsmatrix $\hat{\rho}$ sind durch die Gleichungen

$$\hat{\Lambda}^\top \hat{\rho}^{-1} - \hat{\Lambda}^\top \hat{\rho}^{-1} \mathbf{R} \hat{\rho}^{-1} = \mathbf{O} \quad (6.24)$$

und

$$\text{diag}(\hat{\rho}^{-1} - \hat{\rho}^{-1} \mathbf{R} \hat{\rho}^{-1}) = \mathbf{0} \quad (6.25)$$

gegeben, wobei mit $\text{diag}(\mathbf{X})$ die Diagonale der Matrix \mathbf{X} bezeichnet wird. Die geschätzte Korrelationsmatrix ist dann

$$\hat{\rho} = \hat{\Lambda} \hat{\Lambda}^\top + \hat{\Psi} \quad (6.26)$$

Durch elementare Umformungen der Gleichungen (6.24) und (6.25) erhält man schließlich

$$\text{diag}(\hat{\rho} - \mathbf{R}) = \mathbf{0} \quad , \quad (6.27)$$

d.h. es gilt $\hat{\rho}_{ii} = r_{ii}$ für alle i oder

$$\hat{\psi}_{ii} = r_{ii} - \sum_{j=1}^k \hat{\lambda}_{ij}^2 = 1 - \sum_{j=1}^k \hat{\lambda}_{ij}^2 \quad (6.28)$$

Durch Umformung von Gleichung (6.24) gilt außerdem

$$\hat{\Lambda}^\top = \hat{\mathbf{J}}^{-1} \hat{\Lambda}^\top \hat{\Psi}^{-1} (\mathbf{R} - \hat{\Psi}) \quad (6.29)$$

mit der Diagonalmatrix

$$\hat{\mathbf{J}} = \hat{\Lambda}^\top \hat{\Psi}^{-1} \hat{\Lambda} \quad (6.30)$$

Daraus folgt, dass die Matrix

$$\hat{\mathbf{H}} = \hat{\Lambda}^\top \hat{\Psi}^{-1} (\mathbf{R} - \hat{\Psi}) \hat{\Psi}^{-1} \hat{\Lambda} \quad (6.31)$$

gleich der Matrix $\hat{\mathbf{J}}^2$ ist und daher diagonal ist.

Aus Gleichung (6.29) ist ersichtlich, dass die Elemente der Matrix $\hat{\mathbf{J}}$ bzw. die Quadratwurzel der Elemente von $\hat{\mathbf{H}}$ die Eigenwerte sowie die Zeilen der Matrix $\hat{\Lambda}^\top$ die Eigenvektoren von $\hat{\Psi}^{-1} (\mathbf{R} - \hat{\Psi})$ sind.

Eine **Lösung der Gleichungen** kann iterativ gewonnen werden: Seien $\hat{\Lambda}_{[1]}$ und $\hat{\Psi}_{[1]}$ Approximationen von $\hat{\Lambda}$ und $\hat{\Psi}$, die z.B. durch iterative Kommunalitätsschätzung und Extraktion der Faktoren erhalten werden (siehe voriger Abschnitt). Wenn $\hat{\lambda}_{j[1]}^\top$ die j -te Zeile von $\hat{\Lambda}_{[1]}^\top$ bezeichnet, so berechnet man $\hat{\lambda}_{j[2]}^\top$ durch

$$\hat{\lambda}_{j[2]}^\top = \frac{1}{\sqrt{\hat{h}_j}} \mathbf{u}_j^\top \quad (6.32)$$

mit

$$\mathbf{u}_j^\top = \hat{\boldsymbol{\lambda}}_{j[1]}^\top \hat{\boldsymbol{\Psi}}_{[1]}^{-1} \left(\mathbf{R} - \hat{\boldsymbol{\Psi}}_{[1]} - \sum_{t=1}^{j-1} \hat{\boldsymbol{\lambda}}_{t[2]} \hat{\boldsymbol{\lambda}}_{t[2]}^\top \right) \quad (6.33)$$

und

$$\hat{h}_j = \mathbf{u}_j^\top \hat{\boldsymbol{\Psi}}_{[1]}^{-1} \hat{\boldsymbol{\lambda}}_{j[1]} \quad (6.34)$$

für $j = 1, \dots, k$.

Eine bessere Approximation von $\hat{\psi}_{ii}$ kann durch Gleichung (6.28) erhalten werden. Diese neuen Schätzungen bilden die Matrix $\hat{\boldsymbol{\Psi}}_{[2]}$. Mit den neuen Matrizen $\hat{\boldsymbol{\Lambda}}_{[2]}$ und $\hat{\boldsymbol{\Psi}}_{[2]}$ kann ein neuer Iterationszyklus gestartet werden. Der Prozess wird abgebrochen, wenn sich die Schätzung der Ladungsmatrix nur mehr geringfügig ändert.

Dieses Verfahren konvergiert oft sehr langsam, wobei die exakten Bedingungen für die Konvergenz nicht festgelegt sind.

Ein Nachteil der Maximum-Likelihood-Methode (auch der Hauptfaktorenanalyse) ist, dass die Anzahl k der zu extrahierenden Faktoren vorgegeben werden muss. Dieser Nachteil kann kompensiert werden, indem man einen *Test über die Anzahl der Faktoren* durchführt. Man kann mit einem Test zum Niveau α überprüfen, ob die gewählte Anzahl k der extrahierten gemeinsamen Faktoren zu gering war. Dazu wird die Hypothese

$$H_0 : \hat{\boldsymbol{\rho}} = \mathbf{\Lambda} \mathbf{\Lambda}^\top + \boldsymbol{\Psi} \quad , \quad (6.35)$$

wobei $\mathbf{\Lambda}$ eine $(p \times k)$ -Matrix ist, gegen die Hypothese

$$H_1 : \hat{\boldsymbol{\rho}} = \mathbf{\Lambda}_1 \mathbf{\Lambda}_1^\top \quad (6.36)$$

getestet, wobei $\mathbf{\Lambda}_1$ eine $(p \times p)$ -Matrix ist. Der Plausibilitätsquotiententest zum Testen der H_0 ist $\lambda = l_0/l_1$, wobei l_0 und l_1 die Plausibilitätsfunktionen unter der H_0 bzw. H_1 sind.

Die Statistik $-2 \ln \lambda$ ist asymptotisch χ^2 -verteilt. Die Teststatistik

$$T = \left[n - 1 - \frac{1}{6}(2p + 5) - \frac{2}{3}k \right] \ln \frac{|\hat{\boldsymbol{\rho}}|}{|\mathbf{R}|} \quad (6.37)$$

ist approximativ χ_ν^2 -verteilt mit

$$\nu = \frac{1}{2}[(p - k)^2 - p - k] \quad (6.38)$$

Freiheitsgraden (siehe Lawley und Maxwell, 1971). Da die Anzahl der freien Parameter, $\frac{1}{2}p(p + 1)$, größer sein sollte als die Anzahl der zu schätzenden Parameter, $pk + p - \frac{1}{2}k(k - 1)$, erhält man die Bedingung

$$k \leq \frac{1}{2} \left(1 + 2p - \sqrt{1 + 8p} \right) . \quad (6.39)$$

Die Hypothese H_0 , die besagt, dass k die richtige Anzahl von Faktoren ist, wird also zum Niveau α verworfen, falls

$$T > \chi_{\nu; 1-\alpha}^2 \quad (6.40)$$

gilt.

Die durch die Maximum-Likelihood-Methode (bzw. Hauptfaktorenanalyse) bestimmte Ladungsmatrix $\hat{\mathbf{A}}$ ist i.a. nicht optimal bezüglich der Interpretierbarkeit der k extrahierten orthogonalen Faktoren. Es empfiehlt sich daher, auf $\hat{\mathbf{A}}$ ein orthogonales bzw. schiefwinkeliges Rotationsverfahren anzuwenden (siehe nächster Abschnitt).

6.4 Faktorenrotation

Die Forderungen, dass entweder $\hat{\mathbf{A}}^\top \hat{\mathbf{\Psi}}^{-1} \hat{\mathbf{A}}$ oder $\hat{\mathbf{A}}^\top \hat{\mathbf{A}}$ eine Diagonalmatrix sein soll, bewirken zwar die Eindeutigkeit der Lösung, jedoch sind die k Faktoren in den meisten Fällen nicht interpretierbar. Das Faktorenmuster ändert sich von Stichprobe zu Stichprobe und wird durch Einbeziehung neuer Merkmale entscheidend beeinflusst. Das Koordinatensystem der Faktoren soll jedoch die Merkmale möglichst einfach beschreiben. Das soll durch eine Faktorenrotation erreicht werden. Die Struktur, die man dadurch erlangen möchte, wird *Einfachstruktur* genannt.

Für die Einfachstruktur wurden folgende *fünf Postulate* festgelegt (siehe Thurstone, 1944):

1. Jede Zeile einer Ladungsmatrix soll mindestens eine Null enthalten, d.h. jedes Merkmal wird durch höchstens $k - 1$ Faktoren beschrieben.
2. Jede Spalte einer Ladungsmatrix enthält wenigstens k Nullladungen, d.h. jeder Faktor trägt zur Beschreibung von höchstens $p - k$ der p Merkmale bei.
3. In jedem Spaltenpaar der Ladungsmatrix gibt es mehrere Merkmale, die auf dem einen Faktor eine hohe, auf dem anderen Faktor keine Ladung haben.
4. Wurden mehr als 4 Faktoren extrahiert, so soll jedes beliebige Spaltenpaar für eine große Zahl von Merkmalen in beiden Spalten Null enthalten.
5. Für jedes Spaltenpaar sollten nur wenige Merkmale in beiden Spalten hohe Ladungen haben.

Zum besseren Verständnis wird nun der Begriff der Einfachstruktur geometrisch veranschaulicht.

In den nachfolgenden Abbildungen stellen die Punkte die verschiedenen Merkmale dar.

In Abbildung 6.2 ist ein zweidimensionales Faktorenmuster dargestellt, wie es zustandekommt, wenn die Beziehungen zwischen den Variablen keine besondere Ordnung aufweisen. Diesen Fall nennt man *Zufallskonfiguration*. Hier gibt es keinen Hinweis darauf, wie man die Faktorenachsen legen soll.

In Abbildung 6.3 liegen die Merkmalspunkte deutlich in zwei Gruppen. Durch das rotierte Koordinatensystem \tilde{f}_1, \tilde{f}_2 erhält man eine einfachere Beschreibung der Merkmale als durch das ursprüngliche Koordinatensystem f_1, f_2 . Projiziert man die Merkmalspunkte auf die Achsen, so kann man erkennen, dass die Ladungen auf einer Achse groß, auf der anderen Achse hingegen sehr klein sind. Da die rotierten Achsen

Abbildung 6.2: Zufallskonfiguration

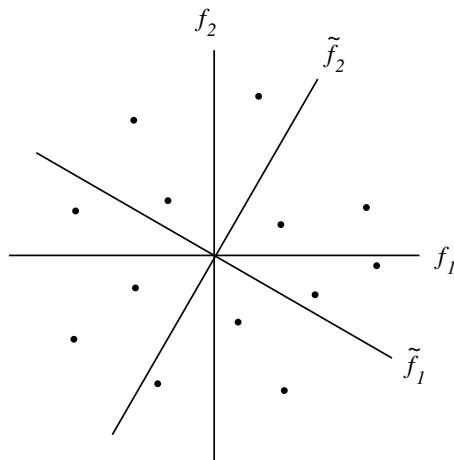
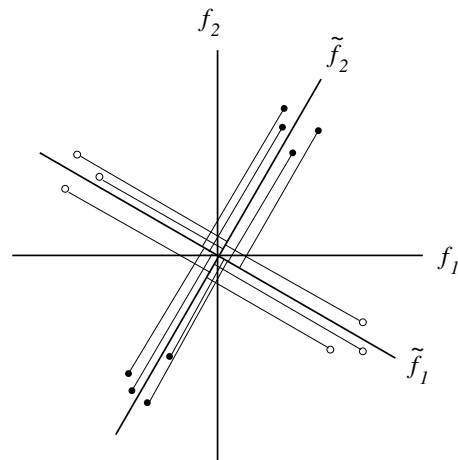


Abbildung 6.3: Orthogonale Einfachstruktur



ebenso wie das ursprüngliche Koordinatensystem aufeinander normal stehen, spricht man von einer *orthogonalen Einfachstruktur*.

Bei realitätsbezogenen Studien gibt es kaum eine so klare Struktur wie in Abbildung 6.3. Die beobachteten Merkmalspunkte sind im Faktorenraum meist in einer Form angeordnet, die ein Mittelding zwischen den beiden Extremen der Zufallskonfiguration und der orthogonalen Einfachstruktur sind. Die Punktwolken liegen meist auf Achsen, die nicht zueinander orthogonal sind. Dieser Fall ist in Abbildung 6.4 veranschaulicht. Man spricht von einer *schiefwinkligen Einfachstruktur*.

Die Projektionen der Punkte auf die Achsen des rotierten Koordinatensystems unterscheiden sich längs einer Achse kaum, auch wenn sie von Punkten des anderen Punktschwarmes stammen. Man kann daher die Punktgruppen schlecht nach der Höhe ihrer Ladungen unterscheiden, wie dies bei der orthogonalen Einfachstruktur der Fall ist.

Abbildung 6.4: Schiefwinklige Einfachstruktur

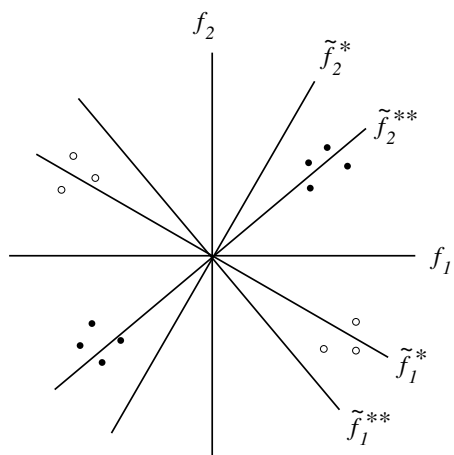
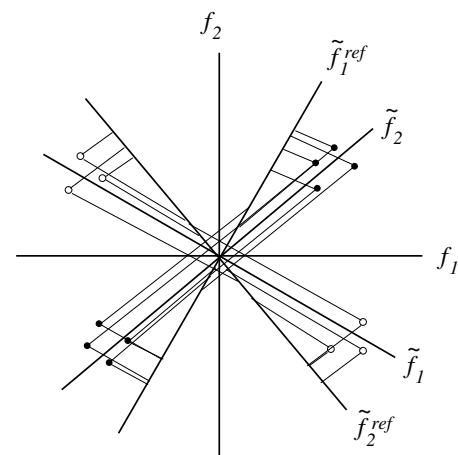


Abbildung 6.5: Referenzvektoren



Dieses Problem wurde mit der Einführung von *Referenzvektoren* \tilde{f}_1^{ref} und \tilde{f}_2^{ref} gelöst (siehe Thurstone, 1961). Durch die Schwerpunkte der Punktwolken werden Achsen gelegt. Die darauf normal stehenden Vektoren (Referenzvektoren) bilden das neue Koordinatensystem. Somit können die Punktgruppen wieder nach der Höhe ihrer Ladungen unterschieden werden (siehe Abbildung 6.5).

Die Faktorenrotation entspricht einer Transformation der Ladungsmatrix. Man unterscheidet zwischen *schiefwinkliger* und *orthogonaler Rotation*.

Bei der orthogonalen Rotation bleibt die Orthogonalität, d.h. die Unkorreliertheit der Faktoren, erhalten. Meist werden jedoch die Forderungen nach Einfachstruktur durch eine schiefwinkelige Rotation, die auch korrelierte Faktoren zulässt, besser erfüllt.

6.4.1 Orthogonale Rotationsverfahren

Ausgegangen wird von der Idee, dass ein Punkt im zweidimensionalen Koordinatensystem dann am einfachsten dargestellt wird, wenn er auf einer der beiden Achsen liegt. Je näher ein Punkt einer Achse kommt, desto kleiner wird das Produkt aus den beiden Koordinaten des Punktes. Da diese Produkte aus den Faktorenladungen positiv oder negativ sein können und die Summe aller Produkte als Kriterium der Einfachheit in Frage kommt, werden die Produkte erst quadriert. Damit ergibt sich folgendes Kriterium der Einfachheit, das ein Minimum annehmen muss, wenn möglichst viele Punkte nahe an den Achsen liegen:

$$\sum_{s < j=1}^k \sum_{i=1}^p \left(\lambda_{is} \lambda_{ij} \right)^2 = Min \quad . \quad (6.41)$$

Die Summe über alle Variablen $i = 1, \dots, p$ und über alle Paare von Faktoren über diesen Ausdruck muss also minimiert werden.

Wenn \mathbf{T} als orthogonale $(k \times k)$ -Transformationsmatrix und $\tilde{\mathbf{\Lambda}}$ als orthogonal rotierte $(p \times k)$ -Ladungsmatrix bezeichnet wird, so wird die ursprüngliche Ladungsmatrix $\mathbf{\Lambda}$ durch die Transformation

$$\tilde{\mathbf{\Lambda}} = \mathbf{\Lambda} \mathbf{T} \quad (6.42)$$

rotiert. Wenn die Transformation orthogonal ist, bleibt die Kommunalität invariant:

$$\kappa_i^2 = \sum_{j=1}^k \lambda_{ij}^2 = \sum_{j=1}^k \tilde{\lambda}_{ij}^2 \quad \text{für } i = 1, \dots, p \quad . \quad (6.43)$$

Auch das Quadrat der Kommunalitäten bleibt durch eine orthogonale Rotation konstant:

$$(\kappa_i^2)^2 = \left(\sum_{j=1}^k \tilde{\lambda}_{ij}^2 \right)^2 = \sum_{j=1}^k \tilde{\lambda}_{ij}^4 + 2 \sum_{s < j=1}^k \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 = konst \quad . \quad (6.44)$$

Auch wenn über alle Variablen p summiert wird, bleibt die erhaltene Größe für orthogonale \mathbf{T} konstant:

$$\sum_{i=1}^p \sum_{j=1}^k \tilde{\lambda}_{ij}^4 + 2 \sum_{i=1}^p \sum_{s < j=1}^k \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 = konst \quad . \quad (6.45)$$

Da die Summe der beiden Ausdrücke konstant ist, wird der eine Summand maximiert, wenn der andere minimiert wird und umgekehrt.

Das bekannte **Quartimax-Kriterium** geht daher von der Maximierung der Funktion

$$QMAX = \sum_{i=1}^p \sum_{j=1}^k \tilde{\lambda}_{ij}^4 \quad (6.46)$$

aus. Nachteil dieser Methode ist, dass ein allgemeiner Faktor zu stark herausgehoben wird; allerdings wird manchmal bewusst ein sogenanntes Einfaktor-Modell gefordert.

Ein anderes und wahrscheinlich das bekannteste orthogonale Rotationskriterium ist das **Varimax-Kriterium**, das von Kaiser (1958) entwickelt wurde. Kaiser betrachtet die Varianz der quadrierten Faktorenladungen für den Faktor j , nämlich

$$s_j^2 = \frac{1}{p} \sum_{i=1}^p \left(\tilde{\lambda}_{ij}^2 \right)^2 - \frac{1}{p^2} \left[\sum_{i=1}^p \tilde{\lambda}_{ij}^2 \right]^2 \quad (6.47)$$

Diese Varianz wird über alle Faktoren summiert. Die so erhaltene Größe muss maximal werden, wenn im gesamten Faktorenmuster die Varianz der quadrierten Faktorenladungen möglichst groß sein soll. Nachteil dieser Größe ist, dass Variablen mit höherer Kommunalität stärker eingehen. Deshalb wird jede Faktorenladung durch die jeweilige Kommunalität dividiert. Das Varimax-Kriterium lautet dann (nach Multiplikation mit p^2)

$$VMAX = p \sum_{j=1}^k \sum_{i=1}^p \left(\frac{\tilde{\lambda}_{ij}}{\kappa_i} \right)^4 - \sum_{j=1}^k \left[\sum_{i=1}^p \left(\frac{\tilde{\lambda}_{ij}}{\kappa_i} \right)^2 \right]^2 \quad (6.48)$$

6.4.2 Schiefwinkelige Rotationsverfahren

Wie bereits in der Einleitung zur Faktorenrotation erwähnt wurde, kann i.a. eine Einfachstruktur durch orthogonale Rotation nicht optimal erreicht werden. Die schiefwinkelige Rotation ist zwar rechenintensiver, aber in Hinblick auf die Interpretierbarkeit der Faktoren ist dieses Verfahren in den meisten Fällen besser geeignet.

Auch hier gibt es eine Reihe von verschiedenen Rotationskriterien, von denen nur zwei vorgestellt werden.

Der gravierende Unterschied zu den orthogonalen Verfahren liegt am Ausgangsmodell

$$\boldsymbol{\rho} = \boldsymbol{\Lambda} \boldsymbol{\Phi} \boldsymbol{\Lambda}^\top + \boldsymbol{\Psi} \quad (6.49)$$

Da die schiefwinkeligen Faktoren miteinander korrelieren, muss auch die Korrelationsmatrix $\boldsymbol{\Phi}$ der Faktoren miteingebracht werden, die bei den orthogonalen Verfahren gleich der Einheitsmatrix ist.

Das zu minimierende **Quartimin-Kriterium** $QMIN$ ist analog zu (6.41) definiert durch

$$QMIN = \sum_{s < j=1}^k \sum_{i=1}^p \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 \quad (6.50)$$

das eine Einfachstruktur der Faktorenladungen bewirken soll.

Die rotierte Ladungsmatrix ist

$$\tilde{\Lambda} = \Lambda \mathbf{T} . \quad (6.51)$$

Durch Einsetzen in das Modell erhält man

$$\begin{aligned} \rho - \Psi &= \Lambda \Phi \Lambda^\top = \tilde{\Lambda} \mathbf{T}^{-1} \Phi (\mathbf{T}^{-1})^\top \tilde{\Lambda}^\top = \tilde{\Lambda} \mathbf{T}^{-1} \text{Cov}(\mathbf{f}) (\mathbf{T}^{-1})^\top \tilde{\Lambda}^\top \\ &= \tilde{\Lambda} \text{Cov}(\mathbf{T}^{-1} \mathbf{f}) \tilde{\Lambda}^\top = \tilde{\Lambda} \text{Cov}(\tilde{\mathbf{f}}) \tilde{\Lambda}^\top . \end{aligned} \quad (6.52)$$

Die rotierten Faktoren $\tilde{\mathbf{f}}$ berechnen sich somit durch

$$\tilde{\mathbf{f}} = \mathbf{T}^{-1} \mathbf{f} . \quad (6.53)$$

Die Transformationsmatrix \mathbf{T} muss somit so konstruiert werden, dass ihre Inverse existiert.

Ein allgemeineres und vielleicht besseres Kriterium als das Quartimin-Kriterium ist das **Oblimin-Kriterium**

$$OBMIN = \sum_{s < j=1}^k \left(\sum_{i=1}^p \tilde{\lambda}_{is}^2 \tilde{\lambda}_{ij}^2 - \frac{\gamma}{p} \sum_{i=1}^p \tilde{\lambda}_{is}^2 \sum_{i=1}^p \tilde{\lambda}_{ij}^2 \right) . \quad (6.54)$$

Durch die Wahl von $\gamma = 0$ erhält man das Quartimin-Kriterium, $\gamma = 1$ liefert das sogenannte *Covarimin-Kriterium*.

Beispiel 6.4.1 *Es werden wiederum die Prüfungsdaten betrachtet. Die durch die Maximum-Likelihood-Methode berechneten Faktoren werden sowohl nach dem Varimax-Kriterium (VMAX) als auch nach dem Quartimin-Kriterium (QMIN) rotiert. Man erhält folgende Ergebnisse:*

Merkmale	VMAX		QMIN	
	$\tilde{\lambda}_{.1}$	$\tilde{\lambda}_{.2}$	$\tilde{\lambda}_{.1}$	$\tilde{\lambda}_{.2}$
ME	0.271	0.675	-0.033	0.752
AG	0.354	0.680	0.078	0.707
LA	0.734	0.513	0.694	0.250
AN	0.742	0.319	0.821	-0.019
ES	0.704	0.285	0.789	-0.041

In Abbildung 6.6 kann man erkennen, dass in diesem Beispiel die orthogonale Rotation zu keiner Verbesserung der Interpretierbarkeit der Faktoren beiträgt. In Abbildung 6.7 hingegen ist erkennbar, dass Elementare Statistik, Analysis und Lineare Algebra durch den Faktor \tilde{f}_1 erklärt werden. Mechanik und Analytische Geometrie werden durch den Faktor \tilde{f}_2 erklärt. Die schiefwinkelige Rotation lässt somit eine Interpretation des Faktors \tilde{f}_1 als algorithmisch abstraktes Denkvermögen und des Faktors \tilde{f}_2 als räumlich abstraktes Denkvermögen zu.

Abbildung 6.6: Varimax-Rotation (orthogonal)

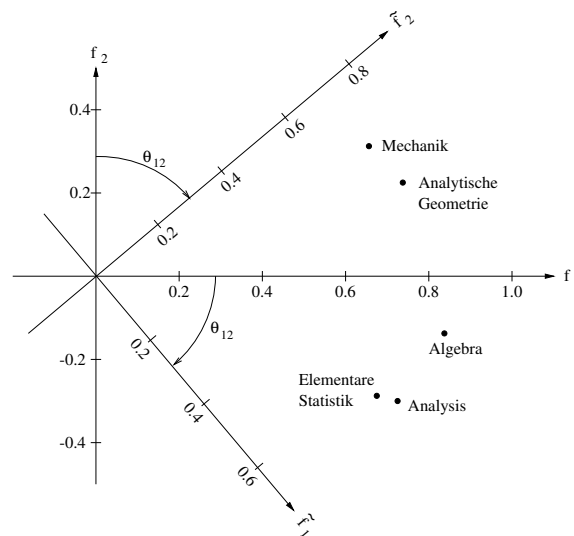
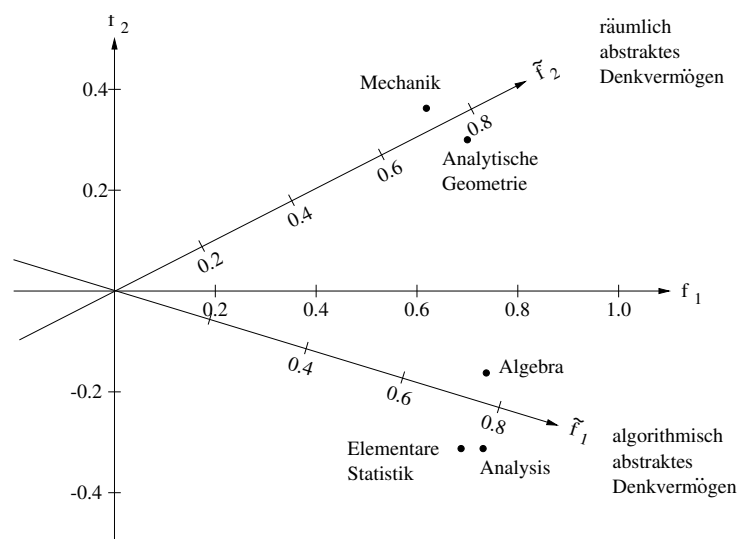


Abbildung 6.7: Quartimin-Rotation (schiefwinkelig)



6.5 Schätzung von Faktorenwerten

Das Hauptinteresse der Faktorenanalyse liegt in der Schätzung der Parameter im Faktorenmodell. Oft möchte man jedoch auch die Werte der gemeinsamen Faktoren (*common factor scores*) ermitteln. Diese Größen können dann für Diagnostik bzw. zur weiteren Analyse (z.B. Clusteranalyse) herangezogen werden.

Die Faktorenwerte sind nicht Schätzungen von unbekannten Parametern im üblichen Sinn. Sie sind Schätzungen von **unbeobachteten** Zufallsvektoren f_j . ($j = 1, \dots, n$). Es gibt verschiedene Möglichkeiten, die Faktorenwerte zu schätzen, zwei der bekanntesten werden nun vorgestellt.

6.5.1 Gewichtete Kleinste-Quadratsummen-Schätzung

Unser Faktorenmodell lautet

$$\mathbf{y} = \mathbf{\Lambda} \mathbf{f} + \mathbf{e} \quad (6.55)$$

mit dem Fehlerterm $\mathbf{e} = (e_1, \dots, e_p)^\top$. Dieses Modell kann auch als Regressionsmodell mit zu schätzenden Regressionsparametern \mathbf{f} aufgefasst werden. Allerdings ist die Forderung nach Homoskedastizität nicht erfüllt, da $\text{Var}(e_i) = \psi_i$ für $i = 1, \dots, p$ nicht notwendigerweise gleich ist (Heteroskedastizität). Man kann das Problem aber leicht in den Griff bekommen, indem man die Gleichung mit $\mathbf{\Psi}^{-1/2}$ multipliziert (gewichtet),

$$\mathbf{\Psi}^{-1/2} \mathbf{y} = \mathbf{\Psi}^{-1/2} \mathbf{\Lambda} \mathbf{f} + \mathbf{\Psi}^{-1/2} \mathbf{e} . \quad (6.56)$$

Für den neuen Fehlerterm $\mathbf{\Psi}^{-1/2} \mathbf{e}$ ist die Gleichheit der Varianzen wieder gegeben, und der Kleinste-Quadratsummenschätzer (*Least Squares*) liefert nun eine unverzerrte Schätzung. Wichtig ist dabei, dass die vorher geschätzten Ladungen und merkmalseigenen Varianzen jetzt als **wahre Werte** behandelt werden. Es ergibt sich

$$\hat{\mathbf{f}} = (\mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{\Lambda})^{-1} \mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{y} \quad (6.57)$$

für die Schätzung der Faktorenwerte (Zufallsvektor). Für eine konkrete Stichprobe sind somit die geschätzten Faktorenwerte der j -ten Beobachtung

$$\hat{\mathbf{f}}_j = (\mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{\Lambda})^{-1} \mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{y}_j . \quad (6.58)$$

Die $(n \times k)$ -Matrix der geschätzten Faktorenwerte ist daher

$$\hat{\mathbf{F}} = \mathbf{Y} \mathbf{\Psi}^{-1} \mathbf{\Lambda} (\mathbf{\Lambda}^\top \mathbf{\Psi}^{-1} \mathbf{\Lambda})^{-1} . \quad (6.59)$$

6.5.2 Regressionsmethode

Wie auch bei der vorigen Methode werden sowohl die Ladungsmatrix $\mathbf{\Lambda}$ als auch die Matrix mit den merkmalseigenen Varianzen $\mathbf{\Psi}$ als bekannt vorausgesetzt. Ausgehend vom Faktorenmodell (6.55) können nun die Faktorenwerte als Regressionsproblem bezüglich der Variablen \mathbf{y} betrachtet werden,

$$\mathbf{f} = \mathbf{B} \mathbf{y} + \boldsymbol{\delta} , \quad (6.60)$$

wobei \mathbf{B} eine $(p \times k)$ -Matrix der Regressionskoeffizienten bezeichnet. Die Regressionskoeffizienten sollen im Sinne der kleinsten Quadrate geschätzt werden. Minimierung von $\boldsymbol{\delta}^\top \boldsymbol{\delta}$ ergibt den Regressionsschätzer

$$\hat{\mathbf{B}} = \mathbf{f} \mathbf{y}^\top (\mathbf{y} \mathbf{y}^\top)^{-1} . \quad (6.61)$$

Da \mathbf{f} unbekannt ist, erhält man eine Schätzung durch

$$\begin{aligned} \hat{\mathbf{f}} &= \hat{\mathbf{B}} \mathbf{y} \\ &= \mathbf{f} \mathbf{y}^\top (\mathbf{y} \mathbf{y}^\top)^{-1} \mathbf{y} \\ &= \mathbf{\Phi} \mathbf{\Lambda}^\top \boldsymbol{\rho}^{-1} \mathbf{y} . \end{aligned} \quad (6.62)$$

ρ ist die Korrelationsmatrix der Variablen. Λ bezeichnet hier sinnvollerweise die rotierte Ladungsmatrix, und daher erhält man mit $\hat{\mathbf{f}}$ auch schon die Faktorenwerte bezüglich des rotierten Koordinatensystems. Für eine konkrete Stichprobe erhält man daher als Schätzung der Faktorenwerte für die j -te Beobachtung ($j = 1, \dots, n$)

$$\hat{\mathbf{f}}_{j.} = \Phi \Lambda^\top \mathbf{R}^{-1} \mathbf{y}_{j.}, \quad (6.63)$$

die gesamte Matrix ist

$$\hat{\mathbf{F}} = \mathbf{Y} \mathbf{R}^{-1} \Lambda \Phi. \quad (6.64)$$

Im Falle von orthogonaler Rotation gilt $\Phi = \mathbf{I}$, womit sich die Formel reduziert auf

$$\hat{\mathbf{F}} = \mathbf{Y} \mathbf{R}^{-1} \Lambda. \quad (\text{orthogonale Faktoren}) \quad (6.65)$$

Es kann auch die Beziehung der Faktorenwerte bezüglich der beiden Schätzmethoden dargestellt werden. Zu diesem Zweck führen wir die Bezeichnung $\hat{\mathbf{f}}^{LS}$ für Kleinst-Quadratsummen-Schätzung und $\hat{\mathbf{f}}^R$ für Regressionsschätzung ein. Man kann leicht zeigen, dass die Identität

$$\Lambda^\top (\Lambda \Phi \Lambda^\top + \Psi)^{-1} = (\mathbf{I} + \Lambda^\top \Psi^{-1} \Lambda \Phi)^{-1} \Lambda^\top \Psi^{-1} \quad (6.66)$$

immer erfüllt ist. Λ soll hier gleich die (schiefwinkelig) rotierte Ladungsmatrix bezeichnen. Es folgt daraus

$$\Lambda^\top \Psi^{-1} \mathbf{y} = (\mathbf{I} + \Lambda^\top \Psi^{-1} \Lambda \Phi) \left[\Lambda^\top (\Lambda \Phi \Lambda^\top + \Psi)^{-1} \mathbf{y} \right]. \quad (6.67)$$

Der Ausdruck in der eckigen Klammer ist aber genau $\Phi^{-1} \hat{\mathbf{f}}^R$, und daher ist

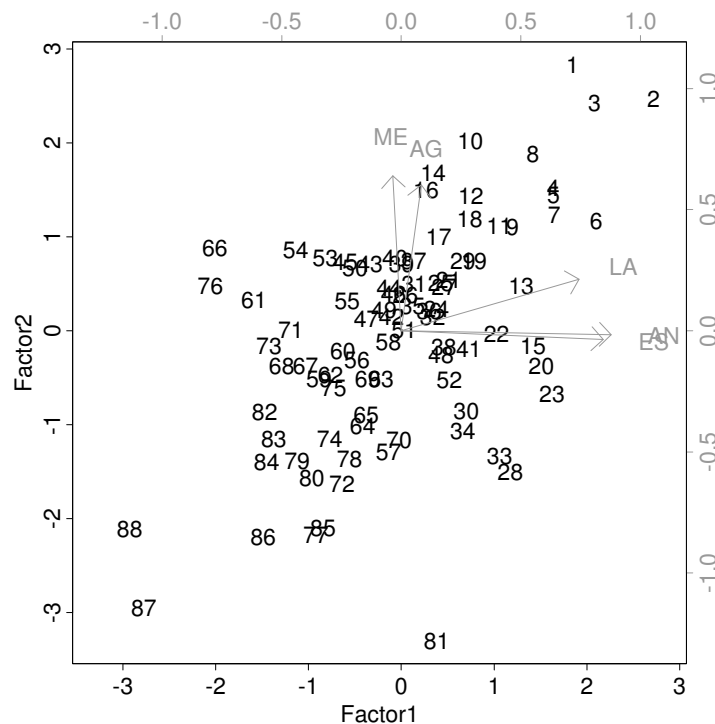
$$\begin{aligned} \hat{\mathbf{f}}^{LS} &= (\Lambda^\top \Psi^{-1} \Lambda)^{-1} (\mathbf{I} + \Lambda^\top \Psi^{-1} \Lambda \Phi) \Phi^{-1} \hat{\mathbf{f}}^R \\ &= \left[(\Phi \Lambda^\top \Psi^{-1} \Lambda)^{-1} + \mathbf{I} \right] \hat{\mathbf{f}}^R \end{aligned} \quad (6.68)$$

Wenn die Elemente der Matrix in runden Klammern nahe 0 sind, werden die beiden Schätzmethoden annähernd das gleiche Ergebnis liefern.

Beispiel 6.5.1 Zur Veranschaulichung der Schätzung von Faktorenwerten wird wieder das Beispiel mit den Prüfungsdaten herangezogen.

Nach erfolgter schiefwinkliger Quartimin-Rotation (siehe letztes Beispiel) können die Faktorenwerte z.B. nach der Gewichteten Kleinst-Quadratsummen-Schätzung berechnet werden. Das Ergebnis ist in Abbildung 6.8 in Form eines Biplots graphisch veranschaulicht. Hohe Werte auf Faktor 1 bedeuten algorithmisch abstraktes Denkvermögen, hohe Werte auf Faktor 2 bedeuten räumlich abstraktes Denkvermögen.

Abbildung 6.8: Schätzung der Faktorenwerte



Literatur

- K. Adel, R. Dutter, H. Filzmoser, and P. Filzmoser. *Tiefenstrukturen der Sprache: Untersuchung regionaler Unterschiede mit statistischen Methoden*. WUV-Universitätsverlag, Wien, 1994.
- D.B. Clarkson and R.I. Jennrich. Quartic Rotation Criteria and Algorithms. *Psychometrika*, 53(2):251–259, 1988.
- W.R. Dillon and M. Goldstein. *Multivariate Analysis: Methods and Applications*. John Wiley & Sons, New York, 1984.
- L. Guttman. Multiple rectilinear prediction and the resolution into components I. *Psychometrika*, 5:75–99, 1940. Cynthia O, Williamsburg, Virginia.
- H.H. Harman. *Modern Factor Analysis*. The University of Chicago Press, Chicago and London, 2nd edition, 1967.
- R.I. Jennrich and P.F. Sampson. Rotation for Simple Loadings. *Psychometrika*, 31(3):313–323, 1966.
- R. Johnson and D. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, London, 4th edition, 1998.
- H.F. Kaiser. The varimax criterion for analytic rotation in factor analysis. *Psychometrika*, 23(3):187–200, 1958.
- D.N. Lawley and A.E. Maxwell. *Factor Analysis as a Statistical Method*. Butterworths, London, 1963.
- D. Revenstorf. *Lehrbuch der Faktorenanalyse*. Verlag W. Kohlhammer, Stuttgart, 1976.

- G.A.F. Seber. *Multivariate observations*. John Wiley & Sons, New York, 1984.
- L.L. Thurstone. Second-order factors. *Psychometrika*, 9:71–100, 1944. Cynthia O, Williamsburg, Virginia.
- L.L. Thurstone. *Multiple factor analysis*. University Press Chicago, Chicago, sixth edition, 1961.
- K. Überla. *Faktorenanalyse*. Springer, Berlin, 1971.

Kapitel 7

Korrelationsanalyse

Bei statistischen Analysen von Daten ist man oft daran interessiert, Zusammenhänge und Abhängigkeiten von Merkmalen festzustellen. Ein Instrumentarium dafür ist die Regressionsanalyse. Möchte man auch das Vorhandensein und die Stärke von Abhängigkeiten messen, so kann die *Korrelationsanalyse* dafür herangezogen werden. Die Korrelation misst den linearen Zusammenhang zwischen Merkmalen.

7.1 Multiple Korrelation

Die *multiple Korrelation* ist ein Maß für die Abhängigkeit eines Merkmals x von einem p -dimensionalen Merkmal $\mathbf{y} = (y_1, \dots, y_p)^\top$. Wir nehmen an, dass sowohl x als auch y_1, \dots, y_p Zufallsvariablen mit gemeinsamer Verteilung sind. Mittelvektor $\boldsymbol{\mu}$ und Kovarianzmatrix $\boldsymbol{\Sigma}$ dieser Verteilung (nicht notwendigerweise Normalverteilung) sind somit

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_x \\ \boldsymbol{\mu}_y \end{pmatrix} \quad \text{bzw.} \quad \boldsymbol{\Sigma} = \begin{pmatrix} \sigma_{xx} & \boldsymbol{\sigma}_{y_x}^\top \\ \boldsymbol{\sigma}_{y_x} & \boldsymbol{\Sigma}_{yy} \end{pmatrix} .$$

Wird nun x durch \mathbf{y} (linear) vorhergesagt, ergibt sich analog zur Regression der Fehler

$$x - a_0 - a_1 y_1 - \dots - a_p y_p .$$

Dieser Fehler ist zufallsabhängig, und daher möchte man die Koeffizienten a_0 und $\mathbf{a} = (a_1, \dots, a_p)^\top$ so wählen, dass der mittlere quadratische Fehler (*mean squared error*)

$$\text{MSE} = E(x - a_0 - \mathbf{a}^\top \mathbf{y})^2$$

minimal wird. Dieser MSE hängt aber von der gemeinsamen Verteilung von x und \mathbf{y} ab, und somit von den Parametern $\boldsymbol{\mu}$ und $\boldsymbol{\Sigma}$.

Satz 7.1.1 Die lineare Vorhersagefunktion $a_0 + \mathbf{a}^\top \mathbf{y}$ mit den Koeffizienten

$$\mathbf{a} = \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\sigma}_{y_x} \quad \text{und} \quad a_0 = \mu_x - \mathbf{a}^\top \boldsymbol{\mu}_y$$

hat minimalen MSE von allen linearen Vorhersagefunktionen von x . Es gilt

$$\text{MSE} = E(x - a_0 - \mathbf{a}^\top \mathbf{y})^2 = E(x - \mu_x - \boldsymbol{\sigma}_{y_x}^\top \boldsymbol{\Sigma}_{yy}^{-1} (\mathbf{y} - \boldsymbol{\mu}_y))^2 = \sigma_{xx} - \boldsymbol{\sigma}_{y_x}^\top \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\sigma}_{y_x} .$$

Außerdem ist

$$a_0 + \mathbf{a}^\top \mathbf{y} = \mu_x + \boldsymbol{\sigma}_{\mathbf{y}x}^\top \boldsymbol{\Sigma}_{\mathbf{y}\mathbf{y}}^{-1} (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{y}})$$

jene lineare Vorhersagefunktion, die maximale Korrelation mit x aufweist, nämlich

$$\text{Corr}(x, a_0 + \mathbf{a}^\top \mathbf{y}) = \sqrt{\frac{\mathbf{a}^\top \boldsymbol{\Sigma}_{\mathbf{y}\mathbf{y}} \mathbf{a}}{\sigma_{xx}}} = \sqrt{\frac{\boldsymbol{\sigma}_{\mathbf{y}x}^\top \boldsymbol{\Sigma}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\sigma}_{\mathbf{y}x}}{\sigma_{xx}}}.$$

Beweis: Siehe z.B. Johnson und Wichern (1998).

Die Korrelation zwischen x und der besten linearen Vorhersagefunktion wird bezeichnet mit *multipler Korrelationskoeffizient* (der Grundgesamtheit)

$$\rho_{x,\mathbf{y}} = + \sqrt{\frac{\boldsymbol{\sigma}_{\mathbf{y}x}^\top \boldsymbol{\Sigma}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\sigma}_{\mathbf{y}x}}{\sigma_{xx}}}.$$

Sein Quadrat $\rho_{x,\mathbf{y}}^2$ wird mit *multipler Bestimmtheitsmaß* (der Grundgesamtheit) bezeichnet, und es gibt an, wie gut das Merkmal x durch die Merkmale $\mathbf{y} = (y_1, \dots, y_p)^\top$ erklärt wird.

Sind anstelle der Kovarianzen Korrelationen gegeben, so kann der multiple Korrelationskoeffizient auch definiert werden als

$$\rho_{x,\mathbf{y}}^2 = \boldsymbol{\rho}_{\mathbf{y}x}^\top \boldsymbol{\rho}_{\mathbf{y}\mathbf{y}}^{-1} \boldsymbol{\rho}_{\mathbf{y}x}.$$

Für konkrete Stichproben kann man die theoretischen Größen durch die entsprechenden Realisierungen schreiben, und man erhält für den multiplen Korrelationskoeffizienten

$$r_{x,\mathbf{y}}^2 = \mathbf{R}_{\mathbf{y}x}^\top \mathbf{R}_{\mathbf{y}\mathbf{y}}^{-1} \mathbf{R}_{\mathbf{y}x}.$$

Ein Test für die Hypothese, dass die multiple Korrelation Null ist, ist gleichwertig dazu, dass alle einfachen Korrelationen Null sind. Man testet also die Hypothese

$$H_0 : \rho_{x,\mathbf{y}} = 0 \quad (= \rho_{xy_1} = \dots = \rho_{xy_p})$$

gegen die Alternativhypothese

$$H_1 : \exists i \in \{1, \dots, p\} \quad \text{mit} \quad \rho_{xy_i} \neq 0$$

zum Niveau α . Wenn n die Anzahl der Objekte in der Stichprobe ist, besitzt die Teststatistik

$$F = \frac{(n-1-p) r_{x,\mathbf{y}}^2}{p (1 - r_{x,\mathbf{y}}^2)}$$

eine $F_{p,n-1-p}$ -Verteilung und die Nullhypothese wird zum Niveau α verworfen, falls

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

gilt.

7.2 Kanonische Korrelation

Bei der kanonischen Korrelation wird die lineare Abhängigkeit zwischen zwei Gruppen von Variablen berechnet. Wie sich herausstellen wird, kann diese Abhängigkeit nicht mehr durch *einen* Korrelationskoeffizienten ausgedrückt werden, sondern es ergibt sich ein Unterraum, der die lineare Abhängigkeit zwischen den Gruppen beschreibt.

Satz 7.2.1 Sei \mathbf{x} eine p -dimensionale und \mathbf{y} eine q -dimensionale Zufallsgröße ($p \leq q$) mit den Erwartungswerten

$$E(\mathbf{x}) = \boldsymbol{\mu}_1 \quad \text{und} \quad E(\mathbf{y}) = \boldsymbol{\mu}_2 .$$

Die Kovarianzmatrizen $\boldsymbol{\Sigma}_{ij}$ mit $i, j = 1, 2$ sind definiert durch

$$\begin{aligned} \boldsymbol{\Sigma}_{11} &= E[(\mathbf{x} - \boldsymbol{\mu}_1)(\mathbf{x} - \boldsymbol{\mu}_1)^\top] , \\ \boldsymbol{\Sigma}_{22} &= E[(\mathbf{y} - \boldsymbol{\mu}_2)(\mathbf{y} - \boldsymbol{\mu}_2)^\top] \quad \text{und} \\ \boldsymbol{\Sigma}_{12} &= \boldsymbol{\Sigma}_{21}^\top = E[(\mathbf{x} - \boldsymbol{\mu}_1)(\mathbf{y} - \boldsymbol{\mu}_2)^\top] \end{aligned}$$

und haben vollen Rang. Wir betrachten die Linearkombinationen $\varphi = \mathbf{a}^\top \mathbf{x}$ und $\eta = \mathbf{b}^\top \mathbf{y}$, wobei \mathbf{a} ein p -dimensionaler und \mathbf{b} ein q -dimensionaler Vektor ist.

Dann ist die einfache Korrelation zwischen φ und η gegeben durch

$$\max_{\mathbf{a}, \mathbf{b}} \text{Corr}(\varphi, \eta) = \rho_1 .$$

Das Maximum wird erreicht durch die Linearkombinationen

$$\varphi_1 = \underbrace{\mathbf{e}_1^\top \boldsymbol{\Sigma}_{11}^{-1/2}}_{\mathbf{a}_1^\top} \mathbf{x} \quad \text{und} \quad \eta_1 = \underbrace{\mathbf{f}_1^\top \boldsymbol{\Sigma}_{22}^{-1/2}}_{\mathbf{b}_1^\top} \mathbf{y} ,$$

die mit **erstes Paar von kanonischen Variablen** bezeichnet werden. ρ_1 heißt **erster kanonische Korrelationskoeffizient**.

Das **k -te Paar von kanonischen Variablen** ($k = 2, 3, \dots, p$) ist gegeben durch

$$\varphi_k = \underbrace{\mathbf{e}_k^\top \boldsymbol{\Sigma}_{11}^{-1/2}}_{\mathbf{a}_k^\top} \mathbf{x} \quad \text{und} \quad \eta_k = \underbrace{\mathbf{f}_k^\top \boldsymbol{\Sigma}_{22}^{-1/2}}_{\mathbf{b}_k^\top} \mathbf{y}$$

und maximiert

$$\text{Corr}(\varphi_k, \eta_k) = \rho_k$$

über alle Linearkombinationen, die unkorreliert mit den vorherigen $1, 2, \dots, k-1$ kanonischen Variablen sind. ρ_k heißt **k -ter kanonische Korrelationskoeffizient**.

$\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$ sind Eigenwerte der Matrix $\boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2}$ und $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$ die zugehörigen Eigenvektoren (Dimension p).

$\rho_1^2, \rho_2^2, \dots, \rho_p^2$ sind auch die p größten Eigenwerte der Matrix $\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1/2}$ mit zugehörigen Eigenvektoren $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_p$ (Dimension q). Jedes \mathbf{f}_i ist proportional zu $\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{e}_i$.

Die kanonischen Variablen haben die Eigenschaften

$$\begin{aligned} \text{Var}(\varphi_k) &= \text{Var}(\eta_k) = 1 \\ \text{Cov}(\varphi_k, \varphi_l) &= \text{Corr}(\varphi_k, \varphi_l) = 0 & k \neq l \\ \text{Cov}(\eta_k, \eta_l) &= \text{Corr}(\eta_k, \eta_l) = 0 & k \neq l \\ \text{Cov}(\varphi_k, \eta_l) &= \text{Corr}(\varphi_k, \eta_l) = 0 & k \neq l \end{aligned}$$

für $k, l = 1, 2, \dots, p$, bzw. mit der Notation $\boldsymbol{\varphi} = (\varphi_1, \dots, \varphi_p)^\top$ und $\boldsymbol{\eta} = (\eta_1, \dots, \eta_p)^\top$ und $\boldsymbol{\rho} = \text{Diag}(\rho_1, \dots, \rho_p)$ gilt somit:

$$\text{Cov} \begin{pmatrix} \boldsymbol{\varphi} \\ \boldsymbol{\eta} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \boldsymbol{\rho} \\ \boldsymbol{\rho} & \mathbf{I} \end{pmatrix}.$$

Beweis: Siehe z.B. Johnson und Wichern (1998).

Bemerkung: Die Koeffizienten \mathbf{a}_i und \mathbf{b}_i , die das Maximierungsproblem von Satz 7.2.1 lösen, können auch auf folgende Weise bestimmt werden:

ρ_i^2 von Satz 7.2.1 ist auch EW von $\boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}$ zum EV $\boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{e}_i = \mathbf{a}_i$,

ρ_i^2 von Satz 7.2.1 ist auch EW von $\boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12}$ zum EV $\boldsymbol{\Sigma}_{22}^{-1/2} \mathbf{f}_i = \mathbf{b}_i$
($i = 1, \dots, p$).

Die kanonischen Korrelationskoeffizienten sind invariant gegenüber linearen Transformationen, was u.a. im folgenden gezeigt wird.

Satz 7.2.2 Ist $\mathbf{x}^* = \mathbf{U}^\top \mathbf{x} + \mathbf{u}$ und $\mathbf{y}^* = \mathbf{V}^\top \mathbf{y} + \mathbf{v}$, wobei \mathbf{U} und \mathbf{V} reguläre $(p \times p)$ - bzw. $(q \times q)$ -Matrizen und \mathbf{u} und \mathbf{v} p - bzw. q -dimensionale feste Vektoren sind. Dann gilt:

- (a) Die kanonischen Korrelationskoeffizienten zwischen \mathbf{x}^* und \mathbf{y}^* sind dieselben wie jene zwischen \mathbf{x} und \mathbf{y} .
- (b) Die Linearkombinationen, die die lineare Abhängigkeit zwischen \mathbf{x}^* und \mathbf{y}^* maximieren, sind gegeben durch

$$\mathbf{a}_i^* = \mathbf{U}^{-1} \mathbf{a}_i \quad \text{und} \quad \mathbf{b}_i^* = \mathbf{V}^{-1} \mathbf{b}_i \quad \text{für } i = 1, \dots, p,$$

wobei \mathbf{a}_i und \mathbf{b}_i die entsprechenden Linearkombinationen für \mathbf{x} und \mathbf{y} sind.

Beweis: Berechnet man statt den Matrizen $\boldsymbol{\Sigma}_{ij}$ für \mathbf{x} und \mathbf{y} die Matrizen $\boldsymbol{\Sigma}_{ij}^*$ für \mathbf{x}^* und \mathbf{y}^* ($i = 1, 2$), so erhält man für $\mathbf{M}_1^* = \boldsymbol{\Sigma}_{11}^{*-1} \boldsymbol{\Sigma}_{12}^* \boldsymbol{\Sigma}_{22}^{*-1} \boldsymbol{\Sigma}_{21}^*$:

$$\mathbf{M}_1^* = (\mathbf{U}^\top \boldsymbol{\Sigma}_{11} \mathbf{U})^{-1} \mathbf{U}^\top \boldsymbol{\Sigma}_{12} \mathbf{V} (\mathbf{V}^\top \boldsymbol{\Sigma}_{22} \mathbf{V})^{-1} \mathbf{V}^\top \boldsymbol{\Sigma}_{21} \mathbf{U} = \mathbf{U}^{-1} \mathbf{M}_1 \mathbf{U}$$

mit $\mathbf{M}_1 = \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}$. Anwendung von Satz 1.4.1 mit $\mathbf{C} = \mathbf{U}^{-1}$ bzw. mit $\mathbf{C} = \mathbf{V}^{-1}$ ergibt die Aussage. \square

Anwendung: Wählt man $\mathbf{U} = (\text{diag} \boldsymbol{\Sigma}_{11})^{-1/2}$ und $\mathbf{V} = (\text{diag} \boldsymbol{\Sigma}_{22})^{-1/2}$, so erhält man anstelle von $\boldsymbol{\Sigma}_{ij}$ die Korrelationsmatrizen $\boldsymbol{\rho}_{ij}$. Statt \mathbf{M}_1 erhält man $\mathbf{M}_1^* = \boldsymbol{\rho}_{11}^{-1} \boldsymbol{\rho}_{12} \boldsymbol{\rho}_{22}^{-1} \boldsymbol{\rho}_{21}$.

Sind also nicht die Kovarianzmatrizen sondern die Korrelationsmatrizen $\boldsymbol{\rho}_{ij}$ gegeben, so sind

- (a) die kanonischen Korrelationskoeffizienten bestimmt durch die Wurzeln aus den Eigenwerten von \mathbf{M}_1^* (oder \mathbf{M}_2^*),
- (b) die Linearkombinationen \mathbf{a}_i und \mathbf{b}_i bestimmt durch die Transformation $\mathbf{a}_i = \mathbf{U}\mathbf{a}_i^*$ bzw. $\mathbf{b}_i = \mathbf{V}\mathbf{b}_i^*$, wobei \mathbf{a}_i^* bzw. \mathbf{b}_i^* die EV von \mathbf{M}_1^* bzw. \mathbf{M}_2^* sind. Es müssen allerdings die Varianzen von x_i und y_i bekannt sein.

Beispiel 7.2.1 *Alle bisher formulierten Sätze lassen sich natürlich auf Stichproben übertragen. Um den Unterschied zwischen Grundgesamtheit und Stichprobe deutlich zu machen, schreiben wir \mathbf{S}_{ij} anstelle von Σ_{ij} , \mathbf{R}_{ij} statt ρ_{ij} und r_i statt ρ_i .*

Gegeben seien die Prüfungsdaten aus Tabelle 2.1. Es wird die kanonische Korrelation zwischen den Fächern Mechanik (ME), Analytische Geometrie (AG) und den Fächern Lineare Algebra (LA), Analysis (AN), Elementare Statistik (ES) berechnet. Die Korrelationsmatrizen der Daten sind

$$\mathbf{R}_{11} = \begin{pmatrix} 1 & 0.5534 \\ 0.5534 & 1 \end{pmatrix}, \quad \mathbf{R}_{22} = \begin{pmatrix} 1 & 0.7108 & 0.6647 \\ 0.7108 & 1 & 0.6072 \\ 0.6647 & 0.6072 & 1 \end{pmatrix},$$

$$\mathbf{R}_{12} = \mathbf{R}_{21}^\top = \begin{pmatrix} 0.5468 & 0.4094 & 0.3891 \\ 0.6096 & 0.4851 & 0.4364 \end{pmatrix}.$$

Aus diesen Korrelationsmatrizen können nach Satz 7.2.2 die Matrizen \mathbf{M}_1^ und \mathbf{M}_2^* berechnet werden. Deren Eigenwerte sind 0.4396 und 0.0017. Die kanonischen Korrelationskoeffizienten sind die Wurzeln aus den Eigenwerten, d.h. $r_1 = 0.6631$ und $r_2 = 0.0409$. Möchte man zusätzlich die kanonischen Variablen erhalten, so berechnet man vorerst \mathbf{a}_i^* und \mathbf{b}_i^* für $i = 1, 2$, wobei in die entsprechenden Formeln statt den Kovarianzmatrizen die Korrelationsmatrizen eingesetzt werden. Es ergibt sich*

$$\mathbf{a}_1^* = \begin{pmatrix} 0.4517 \\ 0.6765 \end{pmatrix}, \quad \mathbf{a}_2^* = \begin{pmatrix} 1.1124 \\ -0.9918 \end{pmatrix},$$

$$\mathbf{b}_1^* = \begin{pmatrix} 0.8703 \\ 0.1191 \\ 0.0596 \end{pmatrix}, \quad \mathbf{b}_2^* = \begin{pmatrix} 0.9600 \\ -1.4608 \\ 0.2473 \end{pmatrix}.$$

Um die Vektoren \mathbf{a}_i und \mathbf{b}_i zu bestimmen, müssen die Matrizen

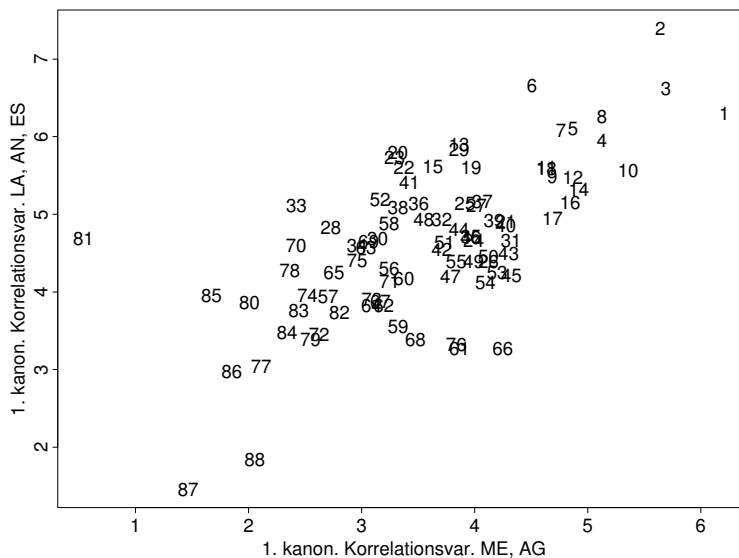
$$\mathbf{U} = (\text{diag} \mathbf{S}_{11})^{-1/2} = \begin{pmatrix} 0.0572 & 0 \\ 0 & 0.0761 \end{pmatrix} \quad \text{und}$$

$$\mathbf{V} = (\text{diag} \mathbf{S}_{22})^{-1/2} = \begin{pmatrix} 0.0941 & 0 & 0 \\ 0 & 0.0674 & 0 \\ 0 & 0 & 0.0580 \end{pmatrix}$$

berechnet werden. Die kanonischen Variablen können nun direkt bestimmt werden, und es ergibt sich

$$\begin{aligned} \varphi_1 &= 0.0258 x_1 + 0.0515 x_2, & \eta_1 &= 0.0819 y_1 + 0.0080 y_2 + 0.0035 y_3, \\ \varphi_2 &= 0.0636 x_1 - 0.0754 x_2, & \eta_2 &= 0.0904 y_1 - 0.0984 y_2 + 0.0143 y_3. \end{aligned}$$

Die beiden ersten kanonischen Variablen sind in Abbildung 7.1 einander gegenübergestellt.

Abbildung 7.1: Erstes Paar von kanonischen Variablen ($r_1 = 0.663$)

In der kanonischen Korrelationsanalyse gibt es eine Reihe von **Tests**, von denen hier nur ausgewählte erwähnt werden, da die Verteilungen der Schätzer sehr kompliziert sind. Mit der Annahme, dass die Daten normalverteilt sind, gilt:

- (a) Ein Plausibilitätsquotiententest für die Hypothese $H_0 : \Sigma_{12} = \mathbf{O}$, d.h. für die Hypothese, dass \mathbf{x} und \mathbf{y} unkorreliert sind, ist durch die Teststatistik

$$\lambda^{2/n} = |\mathbf{I} - \mathbf{S}_{22}^{-1} \mathbf{S}_{21} \mathbf{S}_{11}^{-1} \mathbf{S}_{12}| = \prod_{i=1}^p (1 - r_i^2)$$

gegeben. Diese Teststatistik hat eine $\Lambda(q, n - 1 - p, p)$ -Wilks-Verteilung. n ist der Stichprobenumfang und r_1, \dots, r_p sind die kanonischen Korrelationskoeffizienten.

- (b) Da die Wilks-Verteilung als Produkt von Betaverteilungen relativ kompliziert ist, kann man sie (mit Bartlett's Approximation) für großen Stichprobenumfang n annähern durch

$$- \left[n - \frac{1}{2}(p + q + 3) \right] \ln \prod_{i=1}^p (1 - r_i^2) \sim \chi_{pq}^2 \quad .$$

- (c) Eine ähnliche Teststatistik kann auch für die Hypothese formuliert werden, dass nur s der kanonischen Korrelationskoeffizienten ungleich 0 sind, nämlich

$$- \left[n - \frac{1}{2}(p + q + 3) \right] \ln \prod_{i=s+1}^p (1 - r_i^2) \sim \chi_{(p-s)(q-s)}^2 \quad .$$

Beispiel 7.2.2 Für das Beispiel mit den Prüfungsdaten gilt $n = 88$, $p = 2$ und $q = 3$. Die Statistik zur Hypothese $\rho_1 = \rho_2 = 0$ ist nach (b) gegeben durch

$$-\left(88 - \frac{8}{2}\right) \ln [(1 - 0.6631^2)(1 - 0.0409^2)] \approx 49 .$$

Dieser Wert liegt deutlich über $\chi_{6;0.95}^2 \approx 12.6$, die Hypothese wird daher verworfen. Testet man aber nur die Hypothese $\rho_2 = 0$, so ist die Teststatistik nach (c)

$$-\left(88 - \frac{8}{2}\right) \ln(1 - 0.0409^2) \approx 0.141 .$$

Dies ist beim Test gegen $\chi_{2;0.95}^2 \approx 5.991$ nicht signifikant, weshalb die Hypothese nicht verworfen werden kann.

Literatur

- W.R. Dillon and M. Goldstein. *Multivariate Analysis: Methods and Applications*. John Wiley & Sons, New York, 1984.
- R. Johnson and D. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, London, 4th edition, 1998.
- K.V. Mardia, J.T. Kent, and J.M. Bibby. *Multivariate Analysis*. Acad. Press, London, 1979.
- G.A.F. Seber. *Multivariate observations*. John Wiley & Sons, New York, 1984.

Kapitel 8

Diskriminanzanalyse

8.1 Einleitung

Der Ausdruck *Diskriminanzanalyse* wurde von Fisher (1938) geprägt. Es handelt sich dabei um eine multivariate Methode, die sich einerseits mit der Klassifizierung verschiedener Objektgruppen und andererseits mit der Zuweisung neuer Objekte zu vorher bestimmten Gruppen befasst. Im ersteren Fall wird versucht, die Unterschiede der Objekte, von denen man weiß, dass sie aus zwei oder mehreren Grundgesamtheiten stammen, grafisch oder algebraisch festzuhalten. Man sucht somit nach einer “Diskriminanzfunktion”, die eine möglichst gute Trennung zulässt. Im zweiten Fall möchte man die Objekte in zwei oder mehrere Gruppen einteilen. Das Ziel ist dann, durch festgelegte Regeln *neue* Objekte zu klassifizieren. Obige Fälle stehen oft in unmittelbarem Zusammenhang, denn eine Funktion, die Objekte trennt, kann auch dazu benützt werden, um neue Objekte einzuordnen, bzw. umgekehrt.

8.2 Überlegungen zu Klassifikationsregeln

Die Objekte werden klassifiziert anhand von Messungen von p zugrundeliegenden Zufallsvariablen $\mathbf{X} = (X_1, \dots, X_p)^\top$. Geht man davon aus, dass zwei Gruppen vorhanden sind, so werden die zu messenden Objekte dann eingeteilt in die Klassen π_1 bzw. π_2 . Die Gesamtheit der Werte der ersten Klasse sei die Population der \mathbf{x} -Werte von π_1 , und jene der zweiten Klasse die Population der \mathbf{x} -Werte von π_2 . Die beiden Populationen werden dann beschrieben durch die Wahrscheinlichkeitsverteilungen $f_1(\mathbf{x})$ und $f_2(\mathbf{x})$.

Beispielsweise könnten die beiden Populationen π_1 und π_2 gegeben sein durch Kompositionen von J.S. Bach und L.v. Beethoven. Die zu messenden Variablen \mathbf{X} sind Häufigkeiten verschiedener Akkorde oder Tonfolgen. Es ist dann von Interesse, eine Diskriminanzfunktion zu finden, die beschreibt, inwiefern sich die beiden Komponisten anhand ihrer Werke unterscheiden.

Sei Ω der gesamte Stichprobenraum, also jener Raum, der alle Beobachtungen \mathbf{x} enthält. Jedes Objekt \mathbf{x} muss entweder aus der Population π_1 oder aus π_2 kommen. Weiters sei R_1 jener Raum von Beobachtungen \mathbf{x} , dem wir die Objekte aus π_1

zuordnen, und R_2 jener Raum, dem die restlichen Objekte aus π_2 zugeordnet werden. Ω ist die Vereinigung von R_1 und R_2 .

Bei der Zuordnung kann nun passieren, dass Objekte, die in Wirklichkeit von der Population π_1 kommen, fälschlicherweise als π_2 klassifiziert werden. Sind die Wahrscheinlichkeitsfunktionen $f_1(\mathbf{x})$ und $f_2(\mathbf{x})$ bekannt, so kann diese Wahrscheinlichkeit der falschen Zuordnung als bedingte Wahrscheinlichkeit $P(2|1)$ berechnet werden durch

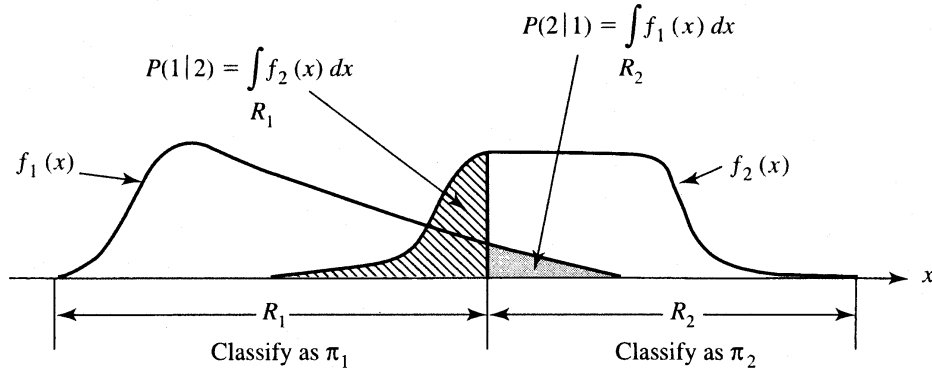
$$P(2|1) = P(\mathbf{X} \in R_2 | \pi_1) = \int_{R_2 = \Omega - R_1} f_1(\mathbf{x}) d\mathbf{x} . \quad (8.1)$$

Umgekehrt kann der Fall eintreten, dass Objekte, die aus der Population π_2 stammen, fälschlicherweise π_1 zugeordnet werden. Die entsprechende Wahrscheinlichkeit ist

$$P(1|2) = P(\mathbf{X} \in R_1 | \pi_2) = \int_{R_1} f_2(\mathbf{x}) d\mathbf{x} . \quad (8.2)$$

Das Integral in (8.1) beschreibt das Volumen der Dichtefunktion $f_1(\mathbf{x})$ über dem Gebiet R_2 , bzw. analog für (8.2).

Abbildung 8.1: Wahrscheinlichkeiten bei Missklassifikation



Sei nun p_1 die Wahrscheinlichkeit, dass die Objekte von π_1 kommen (*a-priori Wahrscheinlichkeit*), und p_2 jene für π_2 , wobei gelten muss $p_1 + p_2 = 1$. Dann können durch Anwenden der Formel für bedingte Wahrscheinlichkeiten folgende Wahrscheinlichkeiten berechnet werden:

$$\begin{aligned} P(\text{Beob. richtig als } \pi_1 \text{ klassifiziert}) &= P(\mathbf{X} \in R_1 | \pi_1) P(\pi_1) = P(1|1) p_1 \\ P(\text{Beob. falsch als } \pi_1 \text{ klassifiziert}) &= P(\mathbf{X} \in R_1 | \pi_2) P(\pi_2) = P(1|2) p_2 \\ P(\text{Beob. richtig als } \pi_2 \text{ klassifiziert}) &= P(\mathbf{X} \in R_2 | \pi_2) P(\pi_2) = P(2|2) p_2 \\ P(\text{Beob. falsch als } \pi_2 \text{ klassifiziert}) &= P(\mathbf{X} \in R_2 | \pi_1) P(\pi_1) = P(2|1) p_1 \end{aligned}$$

Missklassifikation ist oft direkt mit Kosten verbunden. Die Kosten sind natürlich 0, wenn korrekt klassifiziert wurde. Sie betragen $c(1|2)$, wenn eine Beobachtung von π_2 falsch als π_1 klassifiziert wird. Und die Kosten sind $c(2|1)$, wenn Beobachtungen von π_1 fälschlicherweise als π_2 klassifiziert werden. Gemeinsam mit den Wahrscheinlichkeiten für Missklassifikation können nun die *erwarteten Kosten bei Missklassifikation* (EKM) berechnet werden als

$$EKM = c(2|1) P(2|1) p_1 + c(1|2) P(1|2) p_2 . \quad (8.3)$$

Ziel einer Klassifikationsregel ist es, EKM so klein wie möglich zu halten. Demnach könnte eine Klassifikationsregel, die EKM minimiert, folgendermaßen lauten:
Die Menge R_1 ist definiert für Beobachtungen \mathbf{x} , für die gilt:

$$\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \quad (8.4)$$

Die Menge R_2 ist definiert für Beobachtungen \mathbf{x} , für die gilt:

$$\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \quad (8.5)$$

Daraus leiten sich folgende Spezialfälle ab:

(a) $p_1 = p_2$:

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{c(1|2)}{c(2|1)} \quad R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{c(1|2)}{c(2|1)}$$

(b) $c(1|2) = c(2|1)$:

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \frac{p_2}{p_1} \quad R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \frac{p_2}{p_1} \quad (8.6)$$

(c) $p_1 = p_2$ und $c(1|2) = c(2|1)$:

$$R_1 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq 1 \quad R_2 : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < 1$$

Es können aber auch andere Kriterien zur Erstellung einer Klassifikationsregel herangezogen werden. Z.B. könnte man R_1 und R_2 so wählen, dass die *gesamte Wahrscheinlichkeit der Missklassifikation* (GWM) minimal wird, also

$$\begin{aligned} GWM &= P(\text{Missklassifizierung einer Beobachtung aus } \pi_1 \text{ oder } \pi_2) \\ &= p_1 \int_{R_2} f_1(\mathbf{x}) d\mathbf{x} + p_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x} \end{aligned} \quad (8.7)$$

$$= p_1 P(2|1) + p_2 P(1|2) . \quad (8.8)$$

Unter der Voraussetzung, dass die Kosten bei Missklassifikation gleich sind, erkennt man sofort, dass Minimierung von (8.8) äquivalent zur Minimierung von (8.3) ist.

8.3 Der Zweigruppenfall

Wir beschränken uns hier auf multivariat normalverteilte Grundgesamtheiten. Somit sind $f_1(\mathbf{x})$ und $f_2(\mathbf{x})$ Dichtefunktionen multivariater Normalverteilungen mit Mittel $\boldsymbol{\mu}_1$ bzw. $\boldsymbol{\mu}_2$ und Kovarianzmatrizen $\boldsymbol{\Sigma}_1$ bzw. $\boldsymbol{\Sigma}_2$.

8.3.1 Spezialfall $\Sigma_1 = \Sigma_2 = \Sigma$

Die gemeinsame Dichtefunktion der Zufallsvariable $\mathbf{X} = (X_1, \dots, X_p)^\top$ für die Grundgesamtheiten π_1 und π_2 ist gegeben durch

$$f_i(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} |\Sigma|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\} \quad \text{für } i = 1, 2. \quad (8.9)$$

Sind die Parameter $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$ und Σ bekannt, dann ist entsprechend (8.4) der Bereich, der EKM minimiert:

$$\begin{aligned} R_1 : \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right\} \\ \geq \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \\ R_2 : \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right\} \\ < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \end{aligned} \quad (8.10)$$

Anhand dieser Bereiche R_1 und R_2 kann folgende Klassifikationsregel angegeben werden:

Satz 8.3.1 *Seien π_1 und π_2 normalverteilte Populationen mit gleicher Kovarianz. Dann lautet die Klassifikationsregel zur Minimierung von EKM:
Eine Beobachtung \mathbf{x}_0 wird π_1 zugeteilt, wenn*

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1} \mathbf{x}_0 - \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) \geq \ln \left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \right). \quad (8.11)$$

Ansonsten wird \mathbf{x}_0 π_2 zugeordnet.

Beweis: Nachdem alle Größen in (8.10) nichtnegativ sind, ändern sich die Regeln durch Logarithmieren nicht. Weiters gilt:

$$\begin{aligned} & -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \\ &= \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1 - \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \boldsymbol{\mu}_2 - \frac{1}{2} \boldsymbol{\mu}_2^\top \Sigma^{-1} \mathbf{x} + \frac{1}{2} \boldsymbol{\mu}_2^\top \Sigma^{-1} \boldsymbol{\mu}_2 \\ &= \boldsymbol{\mu}_1^\top \Sigma^{-1} \mathbf{x} - \boldsymbol{\mu}_2^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^\top \Sigma^{-1} \boldsymbol{\mu}_2 \end{aligned}$$

Obige Zusammenfassung kann erfolgen, weil die einzelnen Summanden Skalare sind und daher $a = a^\top$ gilt. Erweitert man nun mit $\frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_2$, so ergibt sich weiters:

$$\begin{aligned} & (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_1 - \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_2 + \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} \boldsymbol{\mu}_2 + \frac{1}{2} \boldsymbol{\mu}_2^\top \Sigma^{-1} \boldsymbol{\mu}_2 \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) + \frac{1}{2} \boldsymbol{\mu}_2^\top \Sigma^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2) \end{aligned}$$

Daraus folgt unmittelbar die Aussage. \square

Die in Satz 8.3.1 definierte Diskriminanzfunktion ist *linear* in \mathbf{x} . Meistens sind die Größen $\boldsymbol{\mu}_1$, $\boldsymbol{\mu}_2$ und $\boldsymbol{\Sigma}$ nicht gegeben, und sie müssen aus der Stichprobe geschätzt werden. Seien von der Zufallsvariablen $\mathbf{X} = (X_1, \dots, X_p)^\top$ für die Population π_1 n_1 Stichproben vorhanden und n_2 Messungen für π_2 . Die entsprechenden Datenmatrizen der Dimensionen $(n_1 \times p)$ bzw. $(n_2 \times p)$ seien bezeichnet mit

$$\mathbf{X}_1 = \begin{pmatrix} \mathbf{x}_{11}^\top \\ \mathbf{x}_{12}^\top \\ \vdots \\ \mathbf{x}_{1n_1}^\top \end{pmatrix} \quad \mathbf{X}_2 = \begin{pmatrix} \mathbf{x}_{21}^\top \\ \mathbf{x}_{22}^\top \\ \vdots \\ \mathbf{x}_{2n_2}^\top \end{pmatrix} \quad (8.12)$$

Die arithmetischen Mittelvektoren und empirischen Kovarianzmatrizen sind dann gegeben durch

$$\begin{aligned} \bar{\mathbf{x}}_1 &= \frac{1}{n_1} \sum_{j=1}^{n_1} \mathbf{x}_{1j} & \mathbf{S}_1 &= \frac{1}{n_1 - 1} \sum_{j=1}^{n_1} (\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)(\mathbf{x}_{1j} - \bar{\mathbf{x}}_1)^\top \\ \bar{\mathbf{x}}_2 &= \frac{1}{n_2} \sum_{j=1}^{n_2} \mathbf{x}_{2j} & \mathbf{S}_2 &= \frac{1}{n_2 - 1} \sum_{j=1}^{n_2} (\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)(\mathbf{x}_{2j} - \bar{\mathbf{x}}_2)^\top \end{aligned}$$

Nachdem aber angenommen wurde, dass π_1 und π_2 Populationen mit gleicher Kovarianzmatrix $\boldsymbol{\Sigma}$ sind, werden \mathbf{S}_1 und \mathbf{S}_2 kombiniert zu einer “gepoolten” Kovarianzmatrix

$$\mathbf{S}_{pooled} = \left(\frac{n_1 - 1}{(n_1 - 1) + (n_2 - 1)} \right) \mathbf{S}_1 + \left(\frac{n_2 - 1}{(n_1 - 1) + (n_2 - 1)} \right) \mathbf{S}_2 .$$

\mathbf{S}_{pooled} ist ein unverzerrter Schätzer von $\boldsymbol{\Sigma}$.

Werden nun in (8.11) die Schätzungen eingesetzt, so ergibt sich die Stichproben-Klassifikationsregel:

Satz 8.3.2 *Eine Beobachtung \mathbf{x}_0 wird π_1 zugeteilt, wenn*

$$(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 - \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \geq \ln \left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \right) . \quad (8.13)$$

Ansonsten wird \mathbf{x}_0 π_2 zugeordnet.

Nachdem die Summanden in (8.13) Skalare sind, kann die Regel für den Fall, dass in (8.13) $\left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) = 1$ gilt ($\ln(1) = 0$), auf folgende Weise vereinfacht werden:

$$\hat{y} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x} = \hat{\mathbf{a}}^\top \mathbf{x} \quad (8.14)$$

wird ausgewertet an der Stelle \mathbf{x}_0 und danach verglichen mit der Zahl

$$\hat{m} = \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) = \frac{1}{2} (\bar{y}_1 + \bar{y}_2) ,$$

wobei

$$\bar{y}_1 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \bar{\mathbf{x}}_1 = \hat{\mathbf{a}}^\top \bar{\mathbf{x}}_1$$

und

$$\bar{y}_2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \bar{\mathbf{x}}_2 = \hat{\mathbf{a}}^\top \bar{\mathbf{x}}_2 .$$

Somit reduziert die Regel bei p -dimensionalen Merkmalen die Entscheidung auf eine 1-dimensionale Größe y , die aus entsprechenden Linearkombinationen der Beobachtungen von π_1 und π_2 hervorgeht.

Beispiel 8.3.1 *Dieses Beispiel stammt aus einer Studie von Bouma et al. (1975), die sich mit der Erkennung der Bluterkrankheit beschäftigt. Es gibt verschiedene Typen dieser Krankheit, und die Studie zielt darauf ab, Träger von Hemophilie A zu erkennen. Gemessen wurden die beiden Variablen*

$$\begin{aligned} X_1 &= \log_{10}(\text{AHF Aktivität}) \\ X_2 &= \log_{10}(\text{AHF-ähnliches Antigen}) , \end{aligned}$$

wobei AHF den Anti-Hemophilie-Faktor bezeichnet. Es wurden die Werte von zwei Gruppen von Frauen gemessen, wobei die erste Gruppe mit $n_1 = 30$ Frauen das Hemophilie-Gen nicht trägt und die zweite Gruppe mit $n_2 = 22$ Frauen Träger von Hemophilie A sind. Die Messungen sind in Abbildung 8.2 eingetragen. Es sind außerdem die geschätzten Konturlinien mit Zentrum $\bar{\mathbf{x}}_1$ bzw. $\bar{\mathbf{x}}_2$ eingezeichnet, die 50% bzw. 95% der Beobachtungen enthalten. Die geschätzten Mittel sind

$$\bar{\mathbf{x}}_1 = \begin{pmatrix} -0.0065 \\ -0.0390 \end{pmatrix} \quad \text{und} \quad \bar{\mathbf{x}}_2 = \begin{pmatrix} -0.2483 \\ 0.0262 \end{pmatrix}$$

und die gepoolte Kovarianzmatrix

$$\mathbf{S}_{pooled}^{-1} = \begin{pmatrix} 131.158 & -90.423 \\ -90.423 & 108.147 \end{pmatrix} .$$

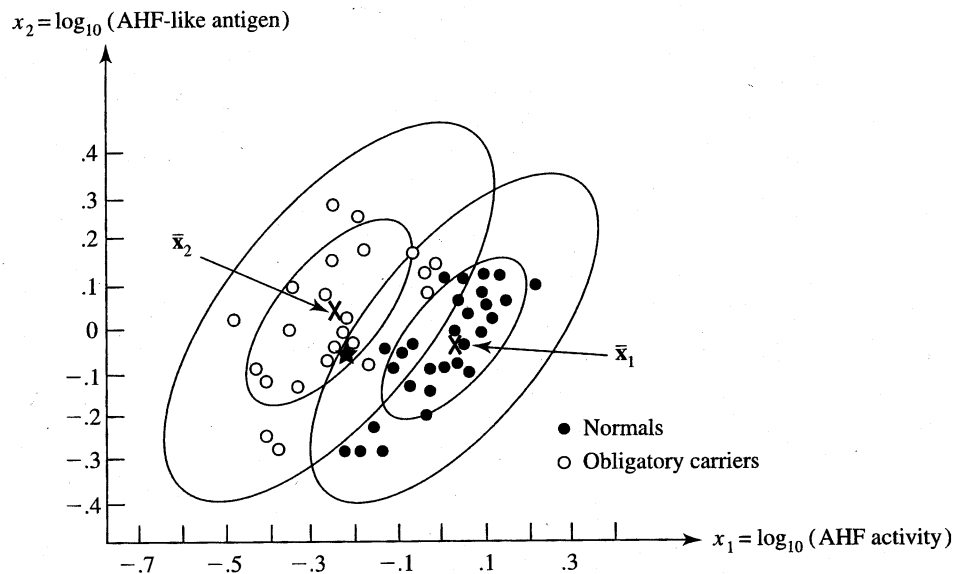
Unter der Annahme, dass die Kosten bei Missklassifikation sowie die a-priori Wahrscheinlichkeiten gleich sind, ergibt sich durch Einsetzen in (8.14)

$$\begin{aligned} \hat{y} &= (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x} = \hat{\mathbf{a}}^\top \mathbf{x} \\ &= (0.2418 \quad -0.0652) \begin{pmatrix} 131.158 & -90.423 \\ -90.423 & 108.147 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \\ &= 37.61x_1 - 28.92x_2 . \end{aligned} \tag{8.15}$$

Weiters sind

$$\begin{aligned} \bar{y}_1 &= \hat{\mathbf{a}}^\top \bar{\mathbf{x}}_1 = (37.61 \quad -28.92) \begin{pmatrix} -0.0065 \\ -0.0390 \end{pmatrix} = 0.88 \\ \bar{y}_2 &= \hat{\mathbf{a}}^\top \bar{\mathbf{x}}_2 = (37.61 \quad -28.92) \begin{pmatrix} -0.2483 \\ 0.0262 \end{pmatrix} = -10.10 \end{aligned}$$

Abbildung 8.2: Plot der Hemophilie Daten



und der Mittelpunkt zwischen diesen Mitteln

$$\hat{m} = \frac{1}{2}(\bar{y}_1 + \bar{y}_2) = \frac{1}{2}(0.88 - 10.10) = -4.61 .$$

Es liegt nun bei einer Frau der Verdacht vor, dass sie Trägerin der Hemophilie A ist. Ihre Werte sind

$$x_1 = -0.210 \quad \text{und} \quad x_2 = -0.044 .$$

Um nun festzustellen, ob sie der Gruppe π_1 (normal) oder der Gruppe π_2 (Träger) zuzuordnen ist, setzen wir in (8.15) die Werte ein und verwenden Klassifikationsregel (8.13):

$$\text{Weise } \mathbf{x}_0 \text{ der Gruppe } \pi_1 \text{ zu, wenn } \hat{y}_0 = \hat{\mathbf{a}}^\top \mathbf{x}_0 \geq \hat{m} = -4.61$$

$$\text{Weise } \mathbf{x}_0 \text{ der Gruppe } \pi_2 \text{ zu, wenn } \hat{y}_0 = \hat{\mathbf{a}}^\top \mathbf{x}_0 < \hat{m} = -4.61$$

mit $\mathbf{x}_0 = (-0.210, -0.044)^\top$. Nachdem

$$\hat{y}_0 = \hat{\mathbf{a}}^\top \mathbf{x}_0 = (37.61 \quad -28.92) \begin{pmatrix} -0.210 \\ -0.044 \end{pmatrix} = -6.62 < -4.61$$

ist, wird die Frau der Gruppe π_2 , also zur Gruppe der Träger der Hemophilie A zugewiesen. Diese Beobachtung ist in Abbildung 8.2 als Stern dargestellt, und man erkennt, dass sie innerhalb des geschätzten 50% Bereichs von π_2 liegt und etwa auf der 95% Konturlinie von π_1 . Die Gruppen sind also nicht klar voneinander getrennt.

Es seien nun die a-priori Wahrscheinlichkeiten p_1 und p_2 für die Gruppenzugehörigkeiten bekannt. Z.B. sei die Frau, die vorhin klassifiziert wurde, eine Cousine mütterlicherseits von einer Trägerin der Hemophilie A. Dann ist nämlich die

genetische Wahrscheinlichkeit, dass auch diese Frau Trägerin ist gleich 0.25. Somit sind $p_1 = 0.75$ (normal) und $p_2 = 0.25$ (Träger). Seien weiters die Kosten bei Missklassifikation gleich (unrealistisch), ergibt sich durch Einsetzen in (8.13)

$$(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 - \frac{1}{2}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1}(\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) = -2.01.$$

Nachdem nun

$$-2.01 < \ln\left(\frac{p_2}{p_1}\right) = \ln\left(\frac{0.25}{0.75}\right) = -1.10$$

ist, wird die Frau wiederum der Gruppe π_2 der Träger von Hemophilie A zugewiesen.

Bemerkung: Die Koeffizienten $\hat{\mathbf{a}} = \mathbf{S}_{pooled}^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$ sind eindeutig bis auf eine multiplikative Konstante, denn jeder Vektor $c\hat{\mathbf{a}}$ mit $c \neq 0$ kann auch als Koeffizientenvektor zur Diskriminierung herangezogen werden. Üblicherweise wird man daher eine Standardisierung vornehmen, z.B.

$$\hat{\mathbf{a}}^* = \frac{\hat{\mathbf{a}}}{\sqrt{\hat{\mathbf{a}}^\top \hat{\mathbf{a}}}}$$

bringt $\hat{\mathbf{a}}^*$ auf Länge 1.

8.3.2 Spezialfall $\Sigma_1 \neq \Sigma_2$

Die bisher formulierten Klassifikationsregeln waren basierend auf (8.4) und haben das Verhältnis der Dichtefunktionen $f_1(\mathbf{x})/f_2(\mathbf{x})$ betrachtet. Beim Fall gleicher Kovarianzmatrizen reduziert sich dieses Verhältnis zu einem relativ einfachen Term, was sich in der Regel (8.11) ausdrückt. Im Fall ungleicher Kovarianzmatrizen (und ungleicher Mittel) wird allerdings das Verhältnis der Dichtefunktionen ein komplizierterer Ausdruck. Wie bei (8.11) kann auch hier der Logarithmus des Verhältnisses betrachtet werden:

$$\begin{aligned} \ln\left(\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})}\right) &= \ln\left(\frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}}\right) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^\top \Sigma_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^\top \Sigma_2^{-1}(\mathbf{x} - \boldsymbol{\mu}_2) \\ &= \frac{1}{2} \ln\left(\frac{|\Sigma_2|}{|\Sigma_1|}\right) - \frac{1}{2} \mathbf{x}^\top (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1^\top \Sigma_1^{-1} - \boldsymbol{\mu}_2^\top \Sigma_2^{-1}) \mathbf{x} \\ &\quad - \frac{1}{2} \boldsymbol{\mu}_1^\top \Sigma_1^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^\top \Sigma_2^{-1} \boldsymbol{\mu}_2 \end{aligned}$$

Damit ergibt sich folgende Klassifikationsregel:

$$\begin{aligned} R_1 : \quad & -\frac{1}{2} \mathbf{x}^\top (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1^\top \Sigma_1^{-1} - \boldsymbol{\mu}_2^\top \Sigma_2^{-1}) \mathbf{x} - k \geq \ln\left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1}\right) \\ R_2 : \quad & -\frac{1}{2} \mathbf{x}^\top (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x} + (\boldsymbol{\mu}_1^\top \Sigma_1^{-1} - \boldsymbol{\mu}_2^\top \Sigma_2^{-1}) \mathbf{x} - k < \ln\left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1}\right) \end{aligned}$$

mit

$$k = \frac{1}{2} \ln \left(\frac{|\Sigma_1|}{|\Sigma_2|} \right) + \frac{1}{2} (\mu_1^\top \Sigma_1^{-1} \mu_1 - \frac{1}{2} \mu_2^\top \Sigma_2^{-1} \mu_2) .$$

Die Konstante k hängt nur mehr von Mittel und Kovarianz der beiden Verteilungen ab. Für den Fall, dass $\Sigma_1 = \Sigma_2$ wird, reduziert sich obige Regel auf (8.11). Es ergibt sich jetzt direkt folgende Aussage:

Satz 8.3.3 *Seien π_1 und π_2 normalverteilte Populationen mit Mittelvektoren μ_1 bzw. μ_2 und Kovarianzmatrizen Σ_1 bzw. Σ_2 . Dann lautet die Klassifikationsregel zur Minimierung von EKM:*

Eine Beobachtung \mathbf{x}_0 wird π_1 zugeteilt, wenn

$$-\frac{1}{2} \mathbf{x}_0^\top (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x}_0 + (\mu_1^\top \Sigma_1^{-1} - \mu_2^\top \Sigma_2^{-1}) \mathbf{x}_0 - k \geq \ln \left(\frac{c(1|2)}{c(2|1)} \frac{p_2}{p_1} \right) . \quad (8.16)$$

Ansonsten wird \mathbf{x}_0 π_2 zugeordnet.

Obige Diskriminanzfunktion ist *quadratisch* in \mathbf{x} . Nachdem die Populationsgrößen μ_1 , μ_2 , Σ_1 und Σ_2 i.a. nicht bekannt sind, werden sie durch die entsprechenden Stichprobengrößen $\bar{\mathbf{x}}_1$, $\bar{\mathbf{x}}_2$, \mathbf{S}_1 und \mathbf{S}_2 geschätzt.

8.3.3 Auswertung der Klassifikation

Obige Klassifikationsregeln gehen von normalverteilten Daten aus. Ist diese Voraussetzung nicht erfüllt, so können die Regeln, die sich durch lineare oder quadratische Diskriminanzfunktionen ausdrücken, stark fehlschlagen. Man könnte daher vorher die Daten so transformieren, dass sie der (multivariaten) Normalverteilung besser entsprechen. Jedenfalls sollte aber überprüft werden, wie gut die Diskriminierung beider Populationen tatsächlich war. Ein Instrumentarium dafür ist die *tatsächliche Fehlerrate* (TFR),

$$TFR = p_1 \int_{\hat{R}_2} f_1(\mathbf{x}) d\mathbf{x} + p_2 \int_{\hat{R}_1} f_2(\mathbf{x}) d\mathbf{x} , \quad (8.17)$$

die sich direkt aus der gesamten Wahrscheinlichkeit der Missklassifikation (GWM) (8.7) herleitet. \hat{R}_1 und \hat{R}_2 sind dabei die aus der Stichprobe ermittelten Bereiche, wenn z.B. Regel (8.13) verwendet wird (Voraussetzung gleiche Kovarianzmatrix).

Man könnte auch einfach den Anteil der falsch klassifizierten Objekte an der Gesamtanzahl der Objekte angeben. Aufgrund einer Stichprobe der Größe n_1 von π_1 und n_2 von π_2 kann dieser Anteil durch Anwenden einer Klassifikationsregel ermittelt werden.

Diese Vorgehensweisen setzen allerdings immer voraus, dass zumindest für einen “Trainingsdatensatz” bekannt ist, welche Objekte welcher Gruppe angehören. Daraus kann die Diskriminanzfunktion ermittelt werden und die Fehlerrate berechnet werden. Neue Objekte können dann entsprechend der vorher ermittelten Regel mit bekannter Fehlerrate klassifiziert werden.

Lachenbruch und Mickey (1968) schlugen folgende “jackknife”-Prozedur zur Schätzung des Fehleranteils vor:

1. Starte mit den Objekten der Gruppe π_1 . Es wird eine Beobachtung dieser Gruppe weggelassen und eine Klassifikationsfunktion mit den restlichen $n_1 - 1$ und n_2 Beobachtungen bestimmt.
2. Die in Schritt 1 weggelassene Beobachtung wird aufgrund der ermittelten Klassifikationsfunktion klassifiziert.
3. Wiederhole Schritt 1 und 2 bis alle Objekte von π_1 klassifiziert sind. \bar{n}_1 sei die Anzahl der falsch klassifizierten Objekte.
4. Wiederhole Schritte 1 bis 3 für die Objekte von π_2 . Es ergibt sich die Anzahl \bar{n}_2 der falsch klassifizierten Objekte von Gruppe π_2 .
5. Nun können die bedingten Wahrscheinlichkeiten für Fehlklassifikation (8.1) und (8.2) geschätzt werden durch

$$\hat{P}(2|1) = \frac{\bar{n}_1}{n_1} \quad \hat{P}(1|2) = \frac{\bar{n}_2}{n_2}$$

und daraus die geschätzte Fehlerrate

$$\frac{\bar{n}_1 + \bar{n}_2}{n_1 + n_2}$$

ermittelt werden.

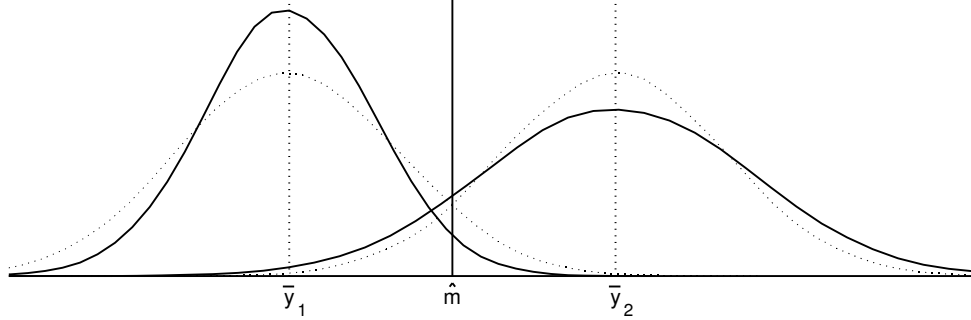
Die praktische Durchführung dieser Prozedur ist möglich bei relativ kleinen Anzahlen n_1 und n_2 bzw. bei Verwendung einer wenig rechenaufwendigen Klassifikationsregel.

Bemerkung: Bei Verwendung der linearen Diskriminanzfunktion aus Satz 8.3.1 ist darauf zu achten, dass die Varianzen der beiden Gruppen tatsächlich als gleich anzunehmen sind. Ist dies nämlich nicht der Fall, kann es passieren, dass die Wahrscheinlichkeiten bei Missklassifikation stark unterschiedlich werden. Dies ist in Abbildung 8.3 dargestellt. Es wurde die lineare Diskriminanzfunktion (8.14) verwendet. Die punktierten Kurven zeigen die Verteilungen der beiden Populationen um \bar{y}_1 und \bar{y}_2 mit der gepoolten Varianz. Anhand des Mittelpunktes \hat{m} zwischen \bar{y}_1 und \bar{y}_2 wird nun für jede Beobachtung entschieden, zu welcher Population sie gezählt wird. Ist nun allerdings die Varianz der Populationen stark unterschiedlich, wie durch die tatsächlichen Verteilungen in der Abbildung angedeutet wird, kann es zu einem starken Ungleichgewicht der tatsächlichen Fehlerwahrscheinlichkeiten kommen. Es würden hier viel eher Werte aus π_2 falsch klassifiziert werden als Werte aus π_1 .

8.3.4 Die Diskriminanzfunktion von Fisher

Fisher (1938) entwickelte eine lineare Diskriminanzfunktion analog zu (8.14), allerdings war die dahinterstehende Idee eine andere. Er versuchte nämlich, multivariate Beobachtungen auf univariate zu transformieren, sodass die beiden transformierten

Abbildung 8.3: Wahrscheinlichkeiten bei Missklassifikation bei Verwendung der linearen Diskriminanzfunktion



Gruppen möglichst stark voneinander separiert sind. Für eine feste Linearkombination der \mathbf{x} -Variablen ergeben sich univariate Werte $y_{11}, y_{12}, \dots, y_{1n_1}$ für die Beobachtungen der ersten Gruppe und Werte $y_{21}, y_{22}, \dots, y_{2n_2}$ für die Beobachtungen der zweiten Gruppe. Die Trennung der beiden Populationen geschieht dann so, dass die arithmetischen Mittel \bar{y}_1 und \bar{y}_2 der univariaten y -Werte möglichst stark voneinander abweichen. Diese Differenz wird ausgedrückt in Einheiten der Standardabweichung, und man erhält daher als Kriterium für die Separation

$$\frac{|\bar{y}_1 - \bar{y}_2|}{s_y}$$

mit der gepoolten Varianz der y -Werte

$$s_y^2 = \frac{\sum_{j=1}^{n_1} (y_{1j} - \bar{y}_1)^2 + \sum_{j=1}^{n_2} (y_{2j} - \bar{y}_2)^2}{n_1 + n_2 - 2}.$$

Es wird nun versucht, eine Linearkombination von \mathbf{x} zu finden, die eine maximale Separation der Stichprobenmittel \bar{y}_1 und \bar{y}_2 zulässt.

Satz 8.3.4 *Die Linearkombination*

$$\hat{y} = \hat{\mathbf{a}}^\top \mathbf{x} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x} \quad (8.18)$$

maximiert das Verhältnis

$$\frac{(\bar{y}_1 - \bar{y}_2)^2}{s_y^2} = \frac{(\hat{\mathbf{a}}^\top \bar{\mathbf{x}}_1 - \hat{\mathbf{a}}^\top \bar{\mathbf{x}}_2)^2}{\hat{\mathbf{a}}^\top \mathbf{S}_{pooled} \hat{\mathbf{a}}} \quad (8.19)$$

über alle möglichen Vektoren $\hat{\mathbf{a}}$. \mathbf{S}_{pooled} ist dabei analog definiert wie im Zweigruppenfall. Das maximale Verhältnis von (8.19) ist

$$D^2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2).$$

Beweis: Mit der Bezeichnung $\mathbf{d} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$ folgt aus der erweiterten Cauchy-Schwarz Ungleichung (Johnson and Wichern, 1998):

$$\max_{\hat{\mathbf{a}}} \frac{(\hat{\mathbf{a}}^\top \mathbf{d})^2}{\hat{\mathbf{a}}^\top \mathbf{S}_{pooled} \hat{\mathbf{a}}} = \mathbf{d}^\top \mathbf{S}_{pooled}^{-1} \mathbf{d} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) = D^2$$

□

Es ergibt sich somit folgende Klassifikationsregel (vgl. auch (8.14)):

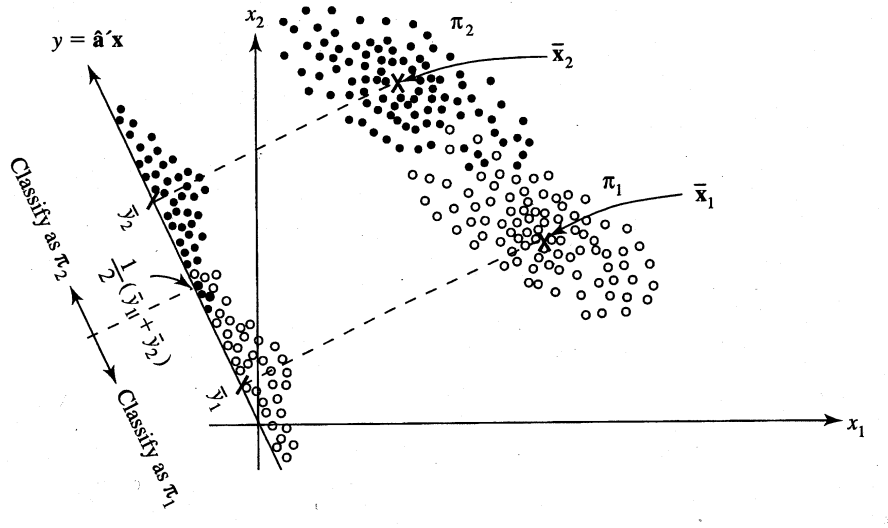
Satz 8.3.5 Eine Beobachtung \mathbf{x}_0 wird π_1 zugeteilt, wenn

$$\hat{y}_0 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 \geq \hat{m} = \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \quad (8.20)$$

Ansonsten wird \mathbf{x}_0 π_2 zugeordnet.

Diese Klassifikationsregel wird für $p = 2$ in Abbildung 8.4 dargestellt. Die Beobachtungen werden auf eine Gerade projiziert. Die Richtung $\hat{\mathbf{a}}$ der Geraden wird so lange variiert, bis die beiden Gruppen maximal separiert sind.

Abbildung 8.4: Schematische Darstellung der Klassifikationsregel von Fisher.



Im Gegensatz zur Regel (8.13) wird bei der Klassifikationsregel von Fisher die Voraussetzung der Normalverteilung beider Populationen nicht benötigt. Allerdings geht man davon aus, dass die Populationen gleiche Kovarianzmatrix haben, nachdem eine gepoolte Schätzung der Kovarianz genommen wird. Man erkennt sofort, dass die lineare Diskriminanzfunktion von Fisher in (8.20) ein spezieller Fall von (8.13) ist. Sind nämlich bei Regel (8.13), die durch Minimierung der erwarteten Kosten bei Missklassifikation hervorgegangen ist, die a-priori Wahrscheinlichkeiten und die erwarteten Kosten bei Missklassifikation gleich, so ergibt sich genau (8.20).

Beispiel 8.3.2 Wir betrachten nochmals die Hemophilie-Daten. Unter der Voraussetzung von gleichen a-priori Wahrscheinlichkeiten und gleichen Kosten bei Missklassifikation ergab sich in (8.15)

$$\hat{y} = \hat{\mathbf{a}}^\top \mathbf{x} = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} \mathbf{x} = 37.61x_1 - 28.92x_2 .$$

Diese Funktion entspricht der linearen Diskriminanzfunktion von Fisher in (8.18) und ermöglicht eine maximale Separation der beiden Gruppen. Die maximale Separation ist gegeben durch

$$\begin{aligned} D^2 &= (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^\top \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2) \\ &= (0.2418 \quad -0.0652) \begin{pmatrix} 131.158 & -90.423 \\ -90.423 & 108.147 \end{pmatrix} \begin{pmatrix} 0.2418 \\ -0.0652 \end{pmatrix} \\ &= 10.98 . \end{aligned}$$

Das heißt also, dass der Wert $D^2 = 10.98$ die maximale Separation der beiden Gruppen angibt, die durch Linearkombination der multivariaten Beobachtungen ermöglicht wird.

8.4 Klassifikation mehrerer Populationen

In diesem Abschnitt sollen nun die Konzepte für den Fall zweier Populationen erweitert werden auf $g > 2$ Gruppen. Diese Erweiterung ist mathematisch plausibel, allerdings ist bisher wenig bekannt über die Eigenschaften der entsprechenden Klassifikationsregeln der Stichprobe. Abweichungen von Normalverteilung der Gruppen oder verschiedene Kovarianzen könnten die guten theoretischen Eigenschaften dann stark verzerren.

8.4.1 Die Methode zur Minimierung der EKM

Sei $f_i(\mathbf{x})$ die Dichte der Beobachtungen von Population π_i mit $i = 1, \dots, g$. Meistens nimmt man an, dass $f_i(\mathbf{x})$ Dichte einer multivariaten Normalverteilung ist, aber dies ist nicht Voraussetzung der nun folgenden Methode.

Die Notation ist angelehnt an den Zweigruppenfall. Demnach bezeichnet p_i die a-priori Wahrscheinlichkeit für die Population π_i ($i = 1, \dots, g$). Weiters sei R_k jener Raum von Beobachtungen \mathbf{x} , dem die Objekte aus π_k zugeordnet werden. $c(k|i)$ sind die Kosten für Missklassifikation, wenn Objekte aus π_i fälschlicherweise π_k zugeordnet werden ($k = 1, \dots, g$), wobei für $k = i$ natürlich $c(i|i) = 0$ gilt. Die Wahrscheinlichkeit dieser falschen Zuordnung ist

$$P(k|i) = P(\mathbf{X} \in R_k | \pi_i) = \int_{R_k} f_i(\mathbf{x}) d\mathbf{x} . \quad (8.21)$$

Die bedingten erwarteten Kosten bei Missklassifikation von \mathbf{x} aus π_1 in π_2, \dots, π_g sind analog zu (8.3)

$$EKM(1) = \sum_{k=2}^g P(k|1) c(k|1) . \quad (8.22)$$

Diese erwarteten Kosten entstehen mit a-priori Wahrscheinlichkeit p_1 . Analog werden $EKM(2), \dots, EKM(g)$ definiert, und man erhält durch Multiplikation mit den a-priori Wahrscheinlichkeiten und Summation aller Beiträge die gesamten erwarteten Kosten bei Missklassifikation (EKM) als

$$EKM = \sum_{i=1}^g p_i EKM(i) = \sum_{i=1}^g p_i \left(\sum_{\substack{k=1 \\ k \neq i}}^g P(k|i) c(k|i) \right). \quad (8.23)$$

Eine optimale Klassifikationsregel sollte nun solche Bereiche R_1, \dots, R_g (disjunkte und vollständige Zerlegung des Stichprobenraumes Ω) ergeben, sodass (8.23) minimiert wird.

Satz 8.4.1 *Die Bereiche zur Minimierung von EKM (8.23) sind dadurch gegeben, dass \mathbf{x} jener Population π_k ($k = 1, \dots, g$) zugeordnet wird, für die der Ausdruck*

$$\sum_{\substack{i=1 \\ k \neq i}}^g p_i f_i(\mathbf{x}) c(k|i) \quad (8.24)$$

minimal wird.

Beweis: siehe Anderson (1984); vgl. auch (8.4)

Nimmt man an, dass die Kosten bei Missklassifikation gleich sind, so würde \mathbf{x} jener Population π_k zugeordnet werden, für die

$$\sum_{\substack{i=1 \\ k \neq i}}^g p_i f_i(\mathbf{x}) \quad (8.25)$$

minimal wird, da in (8.24) die Kosten 1 gesetzt werden können. Minimierung von (8.25) entspricht aber einer Maximierung des weggelassenen Terms $p_k f_k(\mathbf{x})$. Damit vereinfacht sich die Klassifikationsregel zu:

Eine Beobachtung \mathbf{x} wird π_k zugeordnet, wenn

$$p_k f_k(\mathbf{x}) > p_i f_i(\mathbf{x}) \quad \text{für alle } i \neq k. \quad (8.26)$$

Man beachte, dass obige Klassifikationsregeln nur angewandt werden können, wenn die a-priori Wahrscheinlichkeiten, die Missklassifikationskosten und die Dichtefunktionen bekannt sind.

8.4.2 Klassifikation bei Normalverteilung

Als wichtigen Spezialfall betrachten wir multivariat normalverteilte Grundgesamtheiten π_i mit Mittel $\boldsymbol{\mu}_i$ und Kovarianzmatrix $\boldsymbol{\Sigma}_i$,

$$f_i(\mathbf{x}) = (2\pi)^{-\frac{p}{2}} |\boldsymbol{\Sigma}_i|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right\} \quad \text{für } i = 1, \dots, g. \quad (8.27)$$

Verwendet man nun Regel (8.26) bzw. den Logarithmus dieser Regel, so ergibt sich: Eine Beobachtung \mathbf{x} wird π_k zugeordnet, wenn

$$\ln p_k f_k(\mathbf{x}) = \ln p_k - \frac{p}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) = \max_i \ln p_i f_i(\mathbf{x}) . \quad (8.28)$$

Die Konstante $(p/2) \ln(2\pi)$ kann in (8.28) weggelassen werden, da sie gleich ist für alle π_i . Wir definieren daher die *quadratischen Diskriminanzwerte* für die i -te Population als

$$d_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln p_i \quad \text{für } i = 1, \dots, g . \quad (8.29)$$

Somit erhält man folgende Klassifikationsregel:

Eine Beobachtung \mathbf{x} wird π_k zugeordnet, wenn gilt:

$$d_k^Q(\mathbf{x}) \quad \text{ist am größten von } d_1^Q(\mathbf{x}), \dots, d_g^Q(\mathbf{x}) . \quad (8.30)$$

Meistens sind die $\boldsymbol{\mu}_i$ und Σ_i unbekannt und müssen geschätzt werden. Ist ein "Trainingsset" von korrekt klassifizierten Beobachtungen verfügbar, so können mit den Stichprobengrößen n_i der i -ten Population ($i = 1, \dots, g$) die Stichprobenmittel $\bar{\mathbf{x}}_i$ und Stichproben-Kovarianzmatrizen \mathbf{S}_i geschätzt werden. Daraus ergeben sich die geschätzten quadratischen Diskriminanzwerte

$$\hat{d}_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_i)^\top \mathbf{S}_i^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i) + \ln p_i , \quad (8.31)$$

mit deren Hilfe Objekte analog zu (8.30) klassifiziert werden.

Eine Vereinfachung ergibt sich, wenn die Kovarianzmatrizen der Populationen gleich sind, also $\Sigma_i = \Sigma$ für $i = 1, \dots, g$. Dann sind die quadratischen Diskriminanzwerte aus (8.29)

$$d_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} \mathbf{x}^\top \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}_i^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^\top \Sigma^{-1} \boldsymbol{\mu}_i + \ln p_i . \quad (8.32)$$

Nachdem die ersten beiden Terme in (8.32) für alle Populationen gleich sind, können sie weggelassen werden, und man erhält die *linearen Diskriminanzwerte*

$$d_i(\mathbf{x}) = \boldsymbol{\mu}_i^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i^\top \Sigma^{-1} \boldsymbol{\mu}_i + \ln p_i \quad \text{für } i = 1, \dots, g . \quad (8.33)$$

Man erhält eine Schätzung für die linearen Diskriminanzwerte aus der Stichprobe, indem man zunächst eine gepoolte Kovarianzmatrix schätzt durch

$$\mathbf{S}_{pooled} = \frac{1}{n_1 + \dots + n_g - g} \left((n_1 - 1) \mathbf{S}_1 + \dots + (n_g - 1) \mathbf{S}_g \right) . \quad (8.34)$$

Damit erhält man Schätzungen

$$\hat{d}_i(\mathbf{x}) = \bar{\mathbf{x}}_i^\top \mathbf{S}_{pooled}^{-1} \mathbf{x} - \frac{1}{2} \bar{\mathbf{x}}_i^\top \mathbf{S}_{pooled}^{-1} \bar{\mathbf{x}}_i + \ln p_i . \quad (8.35)$$

Man wendet wieder folgende Klassifikationsregel an:

Eine Beobachtung \mathbf{x} wird π_k zugeordnet, wenn gilt:

$$\hat{d}_k(\mathbf{x}) \quad \text{ist am größten von} \quad \hat{d}_1(\mathbf{x}), \dots, \hat{d}_g(\mathbf{x}) . \quad (8.36)$$

Bemerkung: Der Ausdruck (8.33) ist eine lineare Funktion von \mathbf{x} . Man könnte aber eine analoge Regel erhalten, indem man in (8.29) den ersten Term ignoriert, der ja für den Fall gleicher Kovarianzmatrizen für alle Populationen gleich ist. Mit den entsprechenden geschätzten Werten erhält man eine quadratische Distanz

$$D_i^2(\mathbf{x}) = (\mathbf{x} - \bar{\mathbf{x}}_i)^\top \mathbf{S}_{pooled}^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i) , \quad (8.37)$$

und die Klassifikationsregel lautet:

$$\text{Weise } \mathbf{x} \text{ jener Population } \pi_i \text{ zu, für die } -1/2 D_i^2(\mathbf{x}) + \ln p_i \text{ am größten ist.} \quad (8.38)$$

Diese Regel ist analog zur Regel (8.36), sie weist \mathbf{x} jener Population zu, die “am nächsten” ist, wobei das Distanzmaß mit $\ln p_i$ belastet wird. Sind die a-priori Wahrscheinlichkeiten unbekannt, könnte man sie schätzen durch $p_i = 1/g$.

Beispiel 8.4.1 Wir betrachten die bekannten Daten von Ruspini (1970) aus Tabelle 8.1. Die Daten bestehen aus 75 Punkten und wurden ursprünglich von Ruspini verwendet, um Fuzzy-Cluster Methoden zu illustrieren. Die Daten sind dargestellt in Abbildung 8.5, sie enthalten 4 Gruppen, die auch mit Zahlen dargestellt sind. Die Einteilung in die Gruppen erfolgte hier mit k-means Clusterung, und diese entspricht auch einer visuellen Einteilung. Die Anzahlen der Objekte in den einzelnen Gruppen sind

$$n_1 = 17 \quad n_2 = 20 \quad n_3 = 23 \quad n_4 = 15 .$$

Die Gruppenmittel sind mit + eingezeichnet und betragen

$$\bar{\mathbf{x}}_1 = \begin{pmatrix} 98.18 \\ 114.88 \end{pmatrix} \quad \bar{\mathbf{x}}_2 = \begin{pmatrix} 20.15 \\ 64.95 \end{pmatrix} \quad \bar{\mathbf{x}}_3 = \begin{pmatrix} 43.91 \\ 146.04 \end{pmatrix} \quad \bar{\mathbf{x}}_4 = \begin{pmatrix} 68.93 \\ 19.40 \end{pmatrix} .$$

Die Kovarianzmatrizen der einzelnen Gruppen sind

$$\begin{aligned} \mathbf{S}_1 &= \begin{pmatrix} 164.53 & 64.27 \\ 64.27 & 120.36 \end{pmatrix} & \mathbf{S}_2 &= \begin{pmatrix} 92.661 & -5.045 \\ -5.045 & 101.524 \end{pmatrix} \\ \mathbf{S}_3 &= \begin{pmatrix} 91.81 & -23.27 \\ -23.27 & 52.59 \end{pmatrix} & \mathbf{S}_4 &= \begin{pmatrix} 51.2095 & 0.2429 \\ 0.2429 & 52.8286 \end{pmatrix} \end{aligned}$$

Daraus ergibt sich die gepoolte Kovarianzmatrix bzw. deren Inverse als

$$\mathbf{S}_{pooled} = \begin{pmatrix} 100.419 & 5.972 \\ 5.972 & 81.004 \end{pmatrix} \quad \mathbf{S}_{pooled}^{-1} = \begin{pmatrix} 0.0100021 & -0.0007374 \\ -0.0007374 & 0.0123995 \end{pmatrix}$$

Wir möchten nun eine neue Beobachtung $\mathbf{x}_0 = (55, 80)^\top$, die auch in Abbildung 8.5 mit \times eingezeichnet ist, klassifizieren. Wir nehmen dazu an, dass die a-priori

Tabelle 8.1: Ruspini Daten.

x	y	x	y	x	y
4	53	41	150	98	124
5	63	38	145	99	119
10	59	38	143	99	128
9	77	32	143	101	115
13	49	34	141	108	111
13	69	44	156	110	111
12	88	44	149	108	116
15	75	44	143	111	126
18	61	46	142	115	117
19	65	47	149	117	115
22	74	49	152	70	4
27	72	50	142	77	12
28	76	53	144	83	21
24	58	52	152	61	15
27	55	55	155	69	15
28	60	54	124	78	16
30	52	60	136	66	18
31	60	63	139	58	13
32	61	86	132	64	20
36	72	85	115	69	21
28	147	85	96	66	23
32	149	78	94	61	25
35	153	74	96	76	27
33	154	97	122	72	31
38	151	98	116	64	30

Wahrscheinlichkeiten gleich sind, also $p_i = 0.25$. Wir schätzen die linearen Diskriminanzwerte anhand (8.35) und erhalten

$$\hat{d}_1 = 34.42 \quad \hat{d}_2 = 43.08 \quad \hat{d}_3 = 21.98 \quad \hat{d}_4 = 25.81 .$$

Mit der Klassifikationsregel (8.36), die die neue Beobachtung jener Population zuordnet, die am nächsten liegt, wird \mathbf{x}_0 der Population π_2 zugeordnet, weil 43.08 der größte Diskriminanzwert ist.

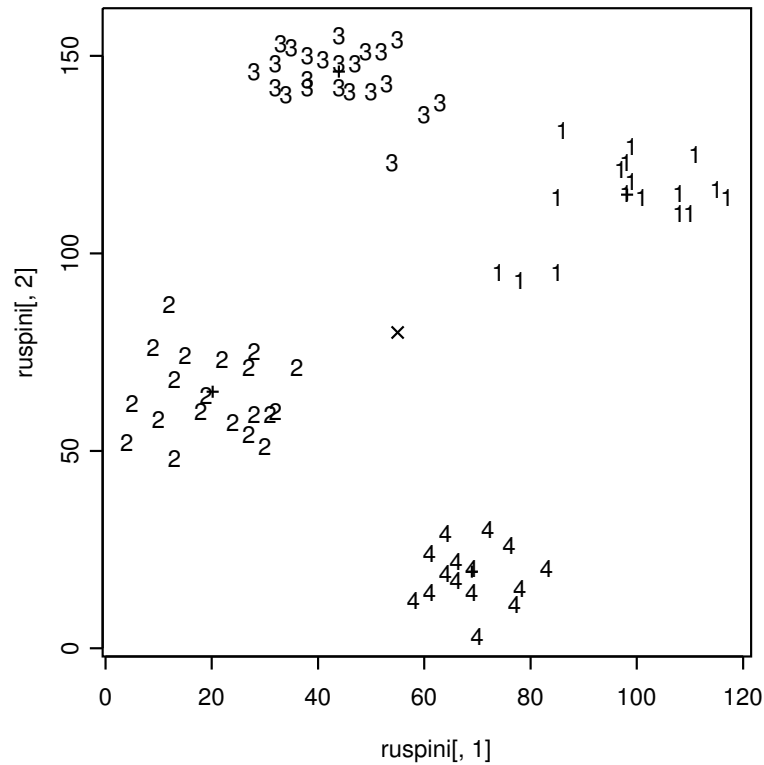
Wenn man nicht davon ausgehen möchte, dass die Kovarianzmatrizen als gleich anzunehmen sind, verwendet man die (geschätzten) quadratischen Diskriminanzwerte aus (8.31) und erhält

$$\hat{d}_1^Q = -13.59 \quad \hat{d}_2^Q = -13.93 \quad \hat{d}_3^Q = -49.06 \quad \hat{d}_4^Q = -42.07 .$$

Der größte Wert ist -13.59, und daher wird \mathbf{x}_0 der Population π_1 zugeordnet. Die Entscheidung fällt aber hier denkbar knapp aus.

Ähnlich zum Zweigruppenfall kann auch hier die “jackknife”-Prozedur von Lachenbruch und Mickey (1968) zur Schätzung des Fehleranteils angewandt werden.

Abbildung 8.5: Ruspini Daten: Erläuterungen siehe Text.



Wenn \bar{n}_i die Anzahl der missklassifizierten Objekte der i -ten Gruppe bezeichnet ($i = 1, \dots, g$), dann ist

$$\frac{\sum_{i=1}^g \bar{n}_i}{\sum_{i=1}^g n_i}$$

die geschätzte Fehlerrate.

Bemerkung: Auch die Methode von Fisher kann auf den Fall mehrerer Gruppen übertragen werden. Die Diskriminanzfunktion ist ein Verhältnis von Variation zwischen den Gruppen zu Variation innerhalb der Gruppen (ähnlich zur Varianzanalyse). Details findet man z.B. in Johnson und Wichern (1998).

8.4.3 Diskriminanzanalyse nach Fisher für den Mehrgruppenfall

Die Diskriminanzanalyse nach Fisher für den Fall $g = 2$ kann erweitert werden auf den Mehrgruppenfall ($g > 2$) siehe Rao (1948). Zu diesem Zweck betrachten wir wieder univariate Projektionen der Form $y = \mathbf{a}^T \mathbf{x}$, mit $\mathbf{a} \in \mathfrak{R}^p$ und $\mathbf{a} \neq \mathbf{0}$, aber diesmal wird eine Projektionsrichtung \mathbf{a} nicht ausreichend sein für die Beschreibung der Lösung.

Sei also $\mathbf{a} \neq \mathbf{0}$ eine gesuchte Projektionsrichtung. Der Erwartungswert für Population π_i (für $i = 1, \dots, g$) der Zufallsvariable $y = \mathbf{a}^T \mathbf{x}$ ist dann

$$\mu_{i,y} = E(y|\mathbf{x} \in \pi_i) = \mathbf{a}^T E(\mathbf{x}|\mathbf{x} \in \pi_i) = \mathbf{a}^T \boldsymbol{\mu}_i$$

und die Varianz ist

$$\sigma_{i,y}^2 = \text{Var}(y|\mathbf{x} \in \pi_i) = \mathbf{a}^T \text{Cov}(\mathbf{x}|\mathbf{x} \in \pi_i) \mathbf{a} = \mathbf{a}^T \boldsymbol{\Sigma}_i \mathbf{a} .$$

Das gesamte gewichtete Mittel der Populationen wird bezeichnet mit $\bar{\boldsymbol{\mu}} = \sum_{i=1}^g p_i \boldsymbol{\mu}_i$, und die entsprechende Projektion ins Univariate mit $\bar{\mu}_y = \mathbf{a}^T \bar{\boldsymbol{\mu}}$.

Wir treffen nun die Annahme, dass die Kovarianzen aller Gruppen gleich sind, also $\boldsymbol{\Sigma}_1 = \dots = \boldsymbol{\Sigma}_g = \boldsymbol{\Sigma}$. Dann gilt auch für obige Varianz

$$\sigma_{1,y}^2 = \dots = \sigma_{g,y}^2 = \sigma_y^2 = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a} .$$

Im Fisher-Zweiggruppenfall wurde ein Ausdruck $(\bar{y}_1 - \bar{y}_2)^2 / s_y^2$ maximiert. In unserer Notation mit Zufallsgrößen entspricht das $(\mu_{1,y} - \mu_{2,y})^2 / \sigma_y^2$. Wir betrachten jetzt zusätzlich a-priori Wahrscheinlichkeiten, und das gewichtete Mittel ist

$$\bar{\mu}_y = p_1 \mu_{1,y} + p_2 \mu_{2,y} = \mathbf{a}^T (p_1 \boldsymbol{\mu}_1 + p_2 \boldsymbol{\mu}_2) .$$

Es ist dann unschwer zu sehen, dass

$$p_1 (\mu_{1,y} - \bar{\mu}_y)^2 + p_2 (\mu_{2,y} - \bar{\mu}_y)^2 = p_1 p_2 (\mu_{1,y} - \mu_{2,y})^2 . \quad (8.39)$$

Letzterer Ausdruck (8.39) sollte also mit der Fisher Regel (bezogen auf die Varianz) maximiert werden, und dieser Ausdruck beschreibt die gewichtete Summe der quadrierten Abstände der Gruppenmittelwerte zum Gesamtmittel.

Die Verallgemeinerung von (8.39) auf den Mehrgruppenfall ist dann unmittelbar einsichtig; es soll nun

$$\frac{\sum_{i=1}^g p_i (\mu_{i,y} - \bar{\mu}_y)^2}{\sigma_y^2} \quad (8.40)$$

maximiert werden. Der Nenner ist laut oben $\sigma_y^2 = \mathbf{a}^T \boldsymbol{\Sigma} \mathbf{a}$. Wenn die Gruppen-Kovarianzen nicht gleich sind, wird die resultierende Klassifikationsregel nicht mehr optimal sein. $\boldsymbol{\Sigma}$ würde dann am besten durch eine gepoolte Version ersetzt werden, also durch

$$\mathbf{W} = \sum_{i=1}^g p_i \boldsymbol{\Sigma}_i .$$

Die Matrix \mathbf{W} beschreibt die *Variation innerhalb der Gruppen*.

Der Zähler von (8.40) kann dargestellt werden als

$$\sum_{i=1}^g p_i (\mu_{i,y} - \bar{\mu}_y)^2 = \sum_{i=1}^g p_i (\mathbf{a}^T (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}}))^2 = \mathbf{a}^T \mathbf{B} \mathbf{a} ,$$

mit der Matrix

$$\mathbf{B} = \sum_{i=1}^g p_i (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})^T ,$$

die die Variation zwischen den Gruppen beschreibt.

Insgesamt kann also das Maximierungsproblem (8.40) ausgedrückt werden als

$$\frac{\mathbf{a}^T \mathbf{B} \mathbf{a}}{\mathbf{a}^T \mathbf{W} \mathbf{a}} \quad \text{für } \mathbf{a} \in \mathfrak{R}^p, \mathbf{a} \neq \mathbf{0} . \quad (8.41)$$

Satz 8.4.2 Die Lösung des Maximierungsproblems (8.41) ist gegeben durch die Eigenvektoren $\mathbf{a}_1, \dots, \mathbf{a}_l$ der Matrix $\mathbf{W}^{-1} \mathbf{B}$, die so skaliert werden, dass $\mathbf{a}_j^T \mathbf{W} \mathbf{a}_j = 1$ für $j = 1, \dots, l$. Die Anzahl l der strikt positiven Eigenwerte der Eigenwertzerlegung von $\mathbf{W}^{-1} \mathbf{B}$ ist dabei $l \leq \min(g-1, p)$.

Beweis: Problem (8.41) ist invariant bezüglich einer Umskalierung von \mathbf{a} . D.h., für jedes $\tilde{\mathbf{a}} = \alpha \mathbf{a}$, mit $\alpha \neq 0$, erhält man das gleiche Maximum. Daher kann man den Nenner so skalieren, dass $\mathbf{a}^T \mathbf{W} \mathbf{a} = 1$ gilt. Dies vereinfacht das Optimierungsproblem (8.41) zu

$$\min_{\mathbf{a}} (-\mathbf{a}^T \mathbf{B} \mathbf{a}) \quad \text{sodass} \quad \mathbf{a}^T \mathbf{W} \mathbf{a} = 1 .$$

Die Minimierung erfolgt mittels Lagrange'schen Ausdruck

$$\phi = -\frac{1}{2} \mathbf{a}^T \mathbf{B} \mathbf{a} + \frac{1}{2} \lambda (\mathbf{a}^T \mathbf{W} \mathbf{a} - 1)$$

mit dem Lagrange-Multiplikator λ . Die Terme $1/2$ sind bequem wenn wir nun die Ableitung bilden:

$$\frac{\partial \phi}{\partial \mathbf{a}} = -\mathbf{B} \mathbf{a} + \lambda \mathbf{W} \mathbf{a}$$

Nullsetzen der Ableitung ergibt

$$\mathbf{B} \mathbf{a} = \mathbf{W} \lambda \mathbf{a} \quad \text{bzw.} \quad \mathbf{W}^{-1} \mathbf{B} \mathbf{a} = \lambda \mathbf{a} .$$

Man erhält also wieder ein Eigenwertproblem, und die Lösungen für \mathbf{a} sind die Eigenvektoren $\mathbf{a}_1, \dots, \mathbf{a}_l$ von $\mathbf{W}^{-1} \mathbf{B}$ zu den Eigenwerten $\lambda_1, \dots, \lambda_l$. Die Eigenwerte sind absteigend sortiert, und dann stellt \mathbf{a}_1 jene Richtung dar, entlang der die Gruppenmittel am besten separiert sind. Somit ist die Reihenfolge der erhaltenen Richtungen wichtig, und mit einer Projektion auf die erste Richtung erhält man – ähnlich wie bei Hauptkomponentenanalyse – die informativste Dimensionsreduktion der gesamten Information, mit dem Unterschied, dass man nun zusätzlich zur multivariaten Information der Beobachtungen auch Klasseninformation hat.

Die Anzahl $l \leq \min(g-1, p)$ erklärt sich unmittelbar aus den Rängen von \mathbf{W} , nämlich maximal p , und \mathbf{B} (maximal $g-1$).

Man bemerke, dass $\mathbf{W}^{-1} \mathbf{B}$ nicht notwendigerweise symmetrisch ist, und somit können Eigenwerte- und vektoren Imaginärteile haben. Dies kann leicht umgangen werden, indem man die symmetrische Matrix \mathbf{B} darstellt als $\mathbf{B} = \mathbf{B}^{1/2} \mathbf{B}^{1/2}$. Mit der Definition $\mathbf{b} = \mathbf{B}^{1/2} \mathbf{a}$, bzw. gleichwertig $\mathbf{a} = \mathbf{B}^{-1/2} \mathbf{b}$ ist obiges Eigenwertproblem $\mathbf{B}^{1/2} \mathbf{W}^{-1} \mathbf{B}^{1/2} \mathbf{b} = \lambda \mathbf{b}$, und man hat somit ein Eigenwertproblem der symmetrischen Matrix $\mathbf{B}^{1/2} \mathbf{W}^{-1} \mathbf{B}^{1/2}$. \square

Man kann nun die *Fisher Diskriminanzfunktionen* definieren als $y_j = \mathbf{a}_j^T \mathbf{x}$, für $j = 1, \dots, l$, die somit die Projektionen der Zufallsgröße \mathbf{x} auf die Richtung \mathbf{a}_j darstellen. Wenn nun konkrete Daten vorliegen, ist eine Visualisierung der beiden ersten

Diskriminanzfunktionen besonders interessant, weil dies jene Projektion der Daten darstellt, in der die Gruppenmittel am besten separiert erscheinen. Man bemerke, dass im Fall von $g = 3$ Gruppen $l \leq 2$ ist, egal ob p groß ist oder nicht.

Schließlich möchte man auch noch eine Klassifikationsregel erhalten. Dazu betrachtet man die *Fisher Diskriminanzwerte*

$$d_i^F(\mathbf{x}) = \sum_{j=1}^l (y_j - \mu_{i,y_j})^2 - 2 \log p_i \quad (8.42)$$

für $i = 1, \dots, g$. Hier ist $\mu_{i,y_j} = \mathbf{a}_j^T \boldsymbol{\mu}_i$, und man hat somit ein Maß für Abweichung von \mathbf{x} zum i -ten Gruppenmittel im Diskriminanzraum, adjustiert mit der a-priori Wahrscheinlichkeit (analog zu früher). Man bemerke, dass hier im Diskriminanzraum das Distanzmaß einfach die Euklidische Distanz ist. Eine neue Beobachtung \mathbf{x} wird dann der Population π_k zugeordnet, wenn $d_k^F(\mathbf{x})$ der kleinste (!) Wert von allen Werten der Gruppen $d_1^F(\mathbf{x}), \dots, d_g^F(\mathbf{x})$ ist.

Durch Anordnen der Eigenvektoren $\mathbf{a}_1, \dots, \mathbf{a}_l$ in den Spalten der Matrix \mathbf{A} kann man die *Fisher Diskriminanzwerte* auch schreiben als

$$d_i^F(\mathbf{x}) = \sum_{j=1}^l (\mathbf{a}_j^T (\mathbf{x} - \boldsymbol{\mu}_i))^2 - 2 \log p_i = (\mathbf{x} - \boldsymbol{\mu}_i)^T \mathbf{A} \mathbf{A}^T (\mathbf{x} - \boldsymbol{\mu}_i) - 2 \log p_i ,$$

was einer (quadratierten) Mahalanobis Distanz im ursprünglichen Raum entspricht.

Literatur

- T.W. Anderson. *An Introduction to Multivariate Statistical Analysis*. John Wiley & Sons, New York, 1984.
- B.N. Bouma, et al. Evaluation of the detection rate of Hemophilia carriers. *Statistical Methods for Clinical Decision Making*, 7(2):339–350, 1975.
- R.A. Fisher. The statistical utilization of multiple measurements. *Annals of Eugenics*, 8:376–386, 1938.
- D.J. Hand. *Discrimination and Classification*. John Wiley & Sons, New York, 1981.
- C.J. Huberty. *Applied Discriminant Analysis*. John Wiley & Sons, New York, 1994.
- R. Johnson and D. Wichern. *Applied Multivariate Statistical Analysis*. Prentice-Hall, London, 4th edition, 1998.
- P.A. Lachenbruch and M.R. Mickey. Estimation of error rates in discriminant analysis. *Technometrics*, 10(1):1–11, 1968.
- C.R. Rao. The utilization to multiple measurements in problems of biological classification. *Journal of the Royal Statistical Society, Series B*, 10:159–203, 1948.
- E.H. Ruspini. Numerical methods for fuzzy clustering. *Inform. Sci.*, 2:319–350, 1970.
- G.A.F. Seber. *Multivariate observations*. John Wiley & Sons, New York, 1984.

Kapitel 9

Projection Pursuit

9.1 Einleitung

9.1.1 Hintergrund

Der Ausdruck *Projection Pursuit* wurde von Friedman und Tukey (1974) eingeführt und bezeichnet eine Methode zur explorativen Analyse multivariater Daten. Mit dieser Methode werden „interessante“ Linearprojektionen der multivariaten Daten auf eine Gerade oder eine Ebene gesucht.

Für Projektionen vom zweidimensionalen in den eindimensionalen Raum ist es möglich, im wesentlichen alle solchen Projektionen zu untersuchen und dann jene herauszusuchen, die von Interesse sind. Im höherdimensionalen Raum ist das nicht mehr durchführbar, für solche Daten muss diese „manuelle“ Vorgangsweise automatisiert werden. Projection Pursuit (PP) ist eine Methode, die interessante Projektionen durch lokale Optimierung über alle Projektionsrichtungen auswählt. Das zugehörige Optimierungskriterium wird mit *Interessantheitsindex* bezeichnet. In der Praxis werden Projektionen in den ein- bzw. zweidimensionalen Raum durchgeführt, da bereits dreidimensionale Punktwolken graphisch schwer darstellbar sind.

Viele klassische multivariate Methoden sind (vom Optimierungskriterium betrachtet) Spezialfälle von PP. Ein Beispiel ist die Hauptkomponentenanalyse, bei der der Interessantheitsindex der Anteil an der Gesamtvarianz ist, der durch die projizierten Daten wiedergegeben wird. Ein anderes Beispiel ist die Quartimax-Methode in der Faktorenanalyse.

9.1.2 Definitionen und Notation

Sei $\mathbf{x} = (x_1, \dots, x_p)^\top$ ein Zufallsvektor im \mathbb{R}^p . Eine Linearprojektion vom \mathbb{R}^p in den \mathbb{R}^k ist eine lineare Abbildung \mathbf{A} , wobei \mathbf{A} eine $(p \times k)$ -Matrix vom Rang k darstellt. Der projizierte Zufallsvektor \mathbf{y} ist dann

$$\mathbf{y} = \mathbf{A}^\top \mathbf{x} \quad , \quad (9.1)$$

wobei $\mathbf{y} \in \mathbb{R}^k$. Wir sprechen von einer *orthonormalen* Projektion, wenn die Spaltenvektoren von \mathbf{A} zueinander orthogonal sind und auf Länge Eins normiert sind, d.h. $\mathbf{A}^\top \mathbf{A} = \mathbf{I}$.

Sei F die Verteilung der p -dimensionalen Zufallsvariable \mathbf{x} , dann ist $\mathbf{y} = \mathbf{A}^\top \mathbf{x}$ eine k -dimensionale Zufallsvariable mit der Verteilung F_A . Für $k = 1$ ist \mathbf{A} nur mehr ein Spaltenvektor, um den Unterschied hervorzuheben schreiben wir daher \mathbf{a} . Die Verteilung von \mathbf{y} wird in diesem Fall mit F_a bezeichnet.

PP sucht nach einer Projektionsmatrix \mathbf{A} , die eine bestimmte „Interessantheitsfunktion“ $Q(F_A)$ maximiert. Da anzunehmen ist, dass mehrere verschiedene Projektionen interessante Ergebnisse liefern, ist man nicht nur am absoluten Extremum, sondern auch an lokalen Extrema interessiert. Q ist ein Funktional im Raum der Verteilungen in \mathbb{R}^k . Um die Notation zu vereinfachen, wird im folgenden anstelle von $Q(F_A)$ einfach $Q(\mathbf{A}^\top \mathbf{x})$ geschrieben.

Für eine feste und bekannte Datenmatrix \mathbf{X} ist $Q(\mathbf{A}^\top \mathbf{X})$ eine Funktion $I(\mathbf{A})$ der Projektionsrichtung \mathbf{A} . Diese Funktion $I(\mathbf{A})$ wird mit *Projektionsindex* bezeichnet.

Es werden nur Funktionen Q berücksichtigt, die *affin invariant* sind, d.h. die durch beliebige reguläre Transformationen

$$Q(\mathbf{A}_1^\top \mathbf{x} + \mathbf{b}_1) = Q(\mathbf{x}) \quad (9.2)$$

die Ausgangsdaten \mathbf{x} unverändert lassen.

Um die späteren Berechnungen zu vereinfachen, wird eine Transformation ausgeübt, welche die Dimension reduziert und die Daten skaliert. Die Kovarianzmatrix wird durch Spektralzerlegung auf die Form

$$\text{Cov}(\mathbf{x}) = \mathbf{\Sigma} = \mathbf{\Gamma} \mathbf{\Lambda} \mathbf{\Gamma}^\top \quad (9.3)$$

gebracht, wobei $\mathbf{\Gamma}$ die orthonormale Matrix der Eigenvektoren und $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$ die $(p \times p)$ -Diagonalmatrix der Eigenwerte von $\mathbf{\Sigma}$ ($\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$) ist. Ist q der Rang von $\mathbf{\Sigma}$, dann sind die Komponenten der standardisierten Zufallsvektoren \mathbf{z} durch die Transformation

$$z_j = \lambda_j^{-1/2} \sum_{i=1}^p \gamma_{ij} (x_i - E(x_i)) \quad \text{für } 1 \leq j \leq q \quad (9.4)$$

gegeben. γ_{ij} sind die Elemente der Matrix $\mathbf{\Gamma}$. Durch diese Transformation ist $E(\mathbf{z}) = \mathbf{0}$ und $\text{Cov}(\mathbf{z}) = \mathbf{I}$.

9.2 Der Projektionsindex von Friedman

9.2.1 Einführung

Der Projektionsindex misst die „Interessantheit“ einer Projektion. Die Frage, ob eine Projektion interessant ist, hängt von der Anwendung ab. Möchte man in den Daten eine Struktur finden, die nicht durch die Korrelationsstruktur erkennbar ist, so erscheint die Normalverteilung einer Projektion am wenigsten von Interesse zu sein. Diese hypothetische Annahme kann auch durch den Zentralen Grenzwertsatz begründet werden, der aussagt, dass die Vereinigung von Verteilungen die Normalverteilung approximiert. Möchte man also Projektionen finden, die eine starke Struktur der Daten aufzeigen, muss der Projektionsindex die Abweichung von der Normalverteilung messen.

9.2.2 Ein- und zweidimensionale Projektion

Friedman (1987) entwickelte einen Projektionsindex, der Abweichungen von der Normalverteilung im Zentrum der Verteilung mehr gewichtet als in den Schwänzen. Die grundlegende Idee ist eine Skalentransformation der projizierten Daten durch die kummulative Verteilungsfunktion. Die transformierte Verteilung wird dann mit der Gleichverteilung, die sich durch Anwendung der selben Transformation auf die Normalverteilung ergibt, verglichen.

Im eindimensionalen Fall von Projection Pursuit sucht man eine Linearkombination

$$y = \mathbf{a}^\top \mathbf{z} \quad , \quad (9.5)$$

sodass die Wahrscheinlichkeitsverteilung $p_a(y)$ von y relativ starke Struktur aufweist. Die oben beschriebene Transformation ist

$$v = 2 \Phi(y) - 1 \quad , \quad (9.6)$$

wobei $\Phi(y)$ die kummulative Standardnormalverteilungsfunktion

$$\Phi(y) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^y e^{-\frac{t^2}{2}} dt \quad (9.7)$$

ist. Durch diese Transformation nimmt v Werte im Intervall $(-1, 1)$ an. Ist y standardnormalverteilt, dann ist v gleichverteilt auf diesem Intervall. Ein Maß für die Abweichung von der Gleichverteilung von v entspricht also einem Maß für die Abweichung von der Normalverteilung von y .

In Abbildung 9.1 ist anhand von Beispielen die Funktionsweise dieses Algorithmus illustriert. Die oberen fünf Bilder gehen von der Standardnormalverteilung aus, die nächsten basieren auf der χ^2 -Verteilung mit 2 Freiheitsgraden, und die unteren fünf Bilder zeigen, wie der Algorithmus eine Bimodalverteilung (Komposition zweier Normalverteilungen) transformiert.

Die Dichtefunktion von v ist

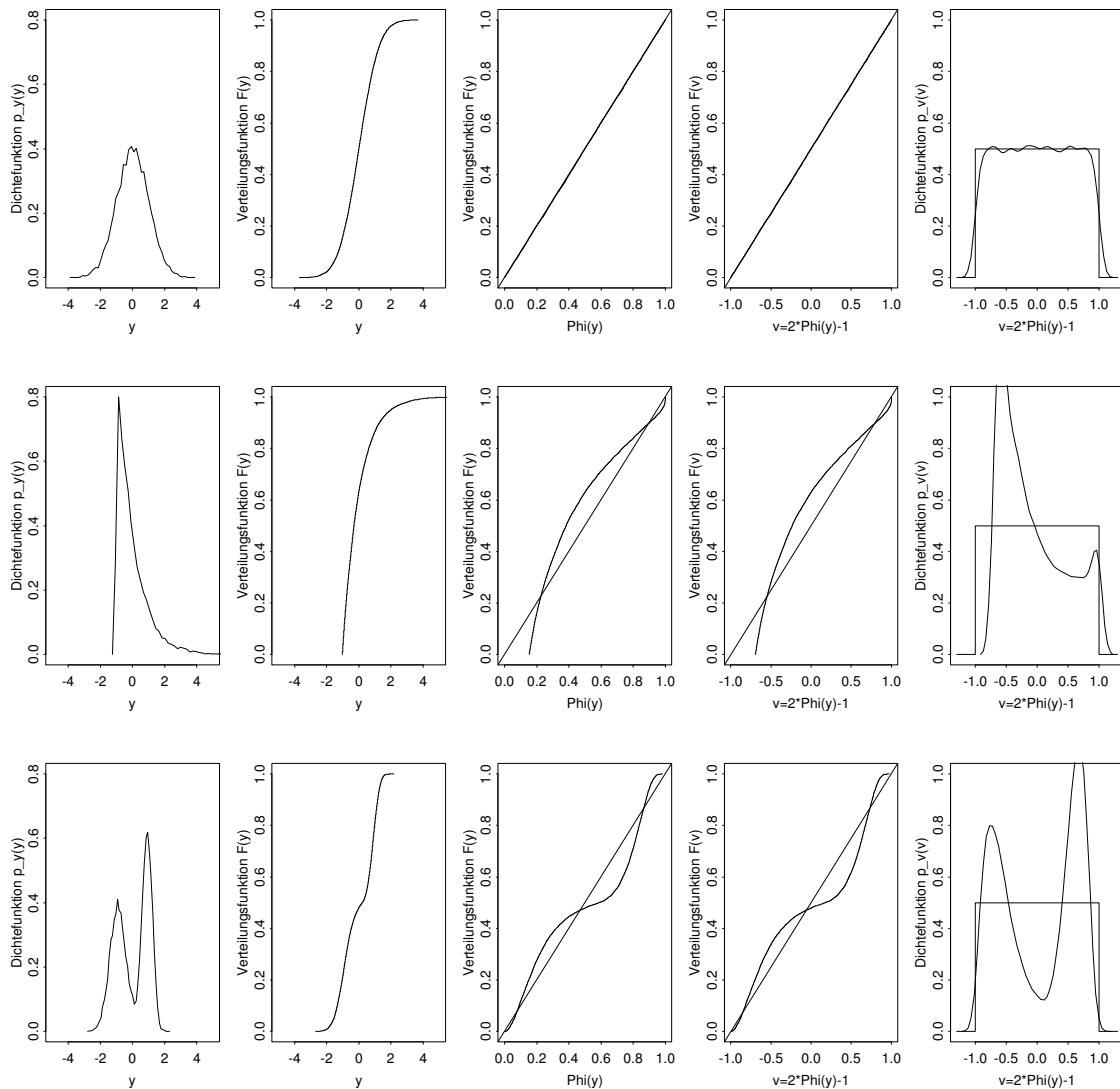
$$p_v(v) = \frac{p_a[\Phi^{-1}(\frac{v+1}{2})]}{2 \phi[\Phi^{-1}(\frac{v+1}{2})]} \quad . \quad (9.8)$$

$\phi(y)$ ist die Dichte der Standardnormalverteilung. Als Maß für die Abweichung von der Gleichverteilung nimmt Friedman das Integral über das Quadrat der Differenz von $p_v(v)$ und der Dichte der Gleichverteilung $p_u(v) = \frac{1}{2}$ über das Intervall $[-1, 1]$:

$$\int_{-1}^1 \left(p_v(v) - \frac{1}{2} \right)^2 dv = \int_{-1}^1 p_v^2(v) dv - \frac{1}{2} \quad . \quad (9.9)$$

Der Projektionsindex $I(\mathbf{a})$ sollte schnell berechenbar sein, sein Wert sollte stetig sein, sodass die ersten Ableitungen bezüglich der Parameter überall existieren. Die Ableitungen sollten ebenfalls schnell berechenbar sein. Diese Forderungen machen es plausibel, dass ein Projektionsindex basierend auf Polynomen am geeignetsten ist.

Abbildung 9.1: Illustration der Funktionsweise des Projektionsindex von Friedman



Die Dichte $p_v(v)$ von v wird durch Legendre-Polynome

$$\int_{-1}^1 p_v^2(v) dv - \frac{1}{2} = \int_{-1}^1 \left(\sum_{j=0}^{\infty} c_j P_j(v) \right) p_v(v) dv - \frac{1}{2} \quad (9.10)$$

approximiert, wobei die Legendre-Polynome durch

$$\begin{aligned} P_0(v) &= 1, \\ P_1(v) &= v, \\ P_j(v) &= \frac{(2j-1)vP_{j-1}(v) - (j-1)P_{j-2}(v)}{j} \quad \text{für } j \geq 2 \end{aligned}$$

definiert werden. Die Koeffizienten c_j sind durch

$$c_j = \frac{2j+1}{2} \int_{-1}^1 P_j(v) p_v(v) dv = \frac{2j+1}{2} E_v[P_j(v)] \quad (9.11)$$

gegeben. Als Approximation für Gleichung (9.9) erhält man somit

$$\int_{-1}^1 p_v^2(v) dv - \frac{1}{2} = \sum_{j=1}^{\infty} \frac{2j+1}{2} E_v^2[P_j(v)] - \frac{1}{2}. \quad (9.12)$$

Für eine Gleichverteilung im Intervall $(-1, 1)$ gilt $E[P_j(v)] = 0$ für $j > 0$.

Den Projektionsindex erhält man, indem die Summe in Gleichung (9.12) bei der Ordnung J abgeschnitten wird, d.h.

$$I(\mathbf{a}) = \sum_{j=1}^J \frac{2j+1}{2} E_v^2[P_j(v)]. \quad (9.13)$$

Für eine effiziente Optimierung ist es wichtig, dass die Ableitungen der Zielfunktion explizit berechenbar sind. Der Projektionsindex muss unter der Bedingung $\mathbf{a}^\top \mathbf{a} = 1$ in Bezug auf alle q Komponenten von \mathbf{a} maximiert werden. Die Ableitungen von $I(\mathbf{a})$ mit $1 \leq k \leq q$ sind

$$\frac{\partial I}{\partial a_k} = \frac{2}{\sqrt{2\pi}} \sum_{j=1}^J (2j+1) E[P_j(v)] E \left[P_j'(v) e^{-\frac{y^2}{2}} (z_k - a_k y) \right] \quad (9.14)$$

mit

$$\begin{aligned} P_1'(v) &= 1, \\ P_j'(v) &= v P_{j-1}'(v) + j P_{j-1}(v) \quad \text{für } j \geq 1. \end{aligned}$$

Für eine Datenmatrix mit n Beobachtungen,

$$\mathbf{Z} = \begin{pmatrix} z_{11} & z_{12} & \dots & z_{1q} \\ z_{21} & z_{22} & \dots & z_{2q} \\ \vdots & \vdots & & \vdots \\ z_{n1} & z_{n2} & \dots & z_{nq} \end{pmatrix},$$

werden die Erwartungswerte in (9.13) durch die entsprechenden Stichprobenmittel geschätzt. Mit den Substitutionen (9.5) und (9.6) erhält man den geschätzten Projektionsindex

$$\hat{I}(\mathbf{a}) = \sum_{j=1}^J \frac{2j+1}{2} \left[\frac{1}{n} \sum_{i=1}^n P_j(2\Phi(\mathbf{a}^\top \mathbf{z}_{i.}) - 1) \right]^2, \quad (9.15)$$

wobei $\mathbf{z}_{i.}$ die i -te Zeile (Spaltenvektor) von \mathbf{Z} bezeichnet.

Der Projektionsindex für die Projektion auf eine Ebene wird analog zum eindimensionalen Fall definiert. Es werden zwei Linearkombinationen

$$y_1 = \mathbf{a}^\top \mathbf{z} , \quad y_2 = \mathbf{b}^\top \mathbf{z} , \quad (9.16)$$

gesucht, sodass die gemeinsame Wahrscheinlichkeitsdichte $p_{ab}(y_1, y_2)$ von y_1 und y_2 relativ starke Struktur aufweist. Da man an nichtlinearen Strukturen interessiert ist, werden die Linearkombinationen y_1 und y_2 als unkorreliert vorausgesetzt. Weil die Daten skaliert wurden, ist diese Bedingung äquivalent dazu, dass \mathbf{a} und \mathbf{b} orthogonal sind, d.h. $\mathbf{a}^\top \mathbf{b} = 0$. Analog zum eindimensionalen Fall sind die weiteren Bedingungen $\mathbf{a}^\top \mathbf{a} = \mathbf{b}^\top \mathbf{b} = 1$.

Die bivariate Normalverteilung wird als die am wenigsten interessierende Verteilung betrachtet. Abweichungen von der Normalverteilung werden im Zentrum stärker gewichtet als in den Schwänzen, man führt daher wieder eine Transformation

$$v_1 = 2 \Phi(y_1) - 1 , \quad v_2 = 2 \Phi(y_2) - 1 \quad (9.17)$$

der projizierten Daten auf das Quadrat $(-1, 1) \times (-1, 1)$ durch. Dadurch sind v_1 und v_2 gleichverteilt, wenn y_1 und y_2 einer bivariaten Normalverteilung folgen. Als Abweichungsmaß kann wieder das Integral über das Quadrat der Differenz zwischen $p_v(v_1, v_2)$ und der Dichte der Gleichverteilung herangezogen werden,

$$\int_{-1}^1 \int_{-1}^1 \left[p_v(v_1, v_2) - \frac{1}{4} \right]^2 dv_1 dv_2 = \int_{-1}^1 \int_{-1}^1 p_v^2(v_1, v_2) dv_1 dv_2 - \frac{1}{4} , \quad (9.18)$$

was einem Maß für die Abweichung von der Gleichverteilung entspricht.

Approximation von $p_v(v_1, v_2)$ mit Legendre-Polynomen liefert

$$\int_{-1}^1 \int_{-1}^1 p_v^2(v_1, v_2) dv_1 dv_2 - \frac{1}{4} = \frac{1}{4} \sum_{j=0}^{\infty} \sum_{l=0}^{\infty} (2j+1)(2l+1) E^2 [P_j(v_1) P_l(v_2)] - \frac{1}{4} . \quad (9.19)$$

Es werden dieselben Legendre-Polynome wie im eindimensionalen Fall verwendet.

Den bivariaten Projektionsindex erhält man durch Abschneiden der Summe der Ordnung J ,

$$\begin{aligned} I(a, b) = & \sum_{j=1}^J \frac{2j+1}{4} E^2[P_j(v_1)] + \sum_{l=1}^J \frac{2l+1}{4} E^2[P_l(v_2)] \\ & + \sum_{j=1}^J \sum_{l=1}^{J-j} \frac{(2j+1)(2l+1)}{4} E^2[P_j(v_1) P_l(v_2)] . \end{aligned}$$

Die Ableitungen des bivariaten Projektionsindex unter den Bedingungen $\mathbf{a}^\top \mathbf{a} = \mathbf{b}^\top \mathbf{b} = 1$ und $\mathbf{a}^\top \mathbf{b} = 0$ mit $1 \leq k \leq q$ und $1 \leq m \leq q$ sind

$$\frac{\partial I}{\partial a_k} = \frac{1}{\sqrt{2\pi}} \sum_{j=1}^J (2j+1) E[P_j(v_1)] E \left[P_j'(v_1) e^{-\frac{y_1^2}{2}} (z_k - a_k y_1 - b_k y_2) \right]$$

$$\begin{aligned}
& + \frac{1}{\sqrt{2\pi}} \sum_{j=1}^J \sum_{l=1}^{J-j} (2j+1)(2l+1) E[P_j(v_1)P_l(v_2)] \\
& E \left[P'_j(v_1)P_l(v_2) e^{-\frac{y_1^2}{2}} (z_k - a_k y_1 - b_k y_2) \right] , \\
\frac{\partial I}{\partial b_m} = & \frac{1}{\sqrt{2\pi}} \sum_{l=1}^J (2l+1) E[P_l(v_2)] E \left[P'_l(v_2) e^{-\frac{y_2^2}{2}} (z_m - a_m y_1 - b_m y_2) \right] \\
& + \frac{1}{\sqrt{2\pi}} \sum_{j=1}^J \sum_{l=1}^{J-j} (2j+1)(2l+1) E[P_j(v_1)P_l(v_2)] \\
& E \left[P_j(v_1)P'_l(v_2) e^{-\frac{y_2^2}{2}} (z_m - a_m y_1 - b_m y_2) \right] .
\end{aligned}$$

y_1 und y_2 wurden in (9.16) definiert, v_1 und v_2 sind durch die Gleichungen (9.17) gegeben. Für eine konkrete Stichprobe werden zur Berechnung des Projektionsindex die Erwartungswerte durch die Stichprobenmittel geschätzt.

Die Berechnungszeit des eindimensionalen Projektionsindex steigt linear mit wachsender Ordnung J der Legendre-Polynome, im zweidimensionalen Fall steigt sie quadratisch.

9.2.3 Entfernung einer Struktur

Der Sinn von PP ist, möglichst viele Projektionen zu finden, die interessante Information enthalten. Wenn man eine solche Projektion gefunden hat, muss die gefundene Struktur entfernt werden, sonst würde immer die selbe Projektion gefunden werden. Es gibt verschiedene Möglichkeiten zum Entfernen von Strukturen (Huber, 1985), ein Weg ist eine rekursive Annäherung wie bei *Projection Pursuit Density Estimation* (Friedman, Stuetzle, und Schroeder, 1984). Die interessante Struktur wird entfernt und der Projektionsindex wird wieder maximiert. Der dabei verwendete Algorithmus ist sehr rechenintensiv.

Friedman (1987) schlägt zum Entfernen einer Struktur einen Algorithmus vor, der schneller berechenbar ist und auch für zweidimensionale Projektionen entwickelt wurde.

Nach der oben aufgestellten Hypothese ist man an Projektionen, deren Dichte der Dichtefunktion der Standardnormalverteilung nahek kommt, nicht interessiert. Die Grundidee für die Entfernung einer Struktur ist daher, durch eine Transformation die Dichte der Projektion in eine Standardnormalverteilung überzuführen. Diese Transformation darf allerdings nur im Unterraum der Projektion wirksam sein und muss alle orthogonalen Richtungen unverändert lassen.

Die Prozedur wird erst für eindimensionale Projektionen beschrieben. Sei \mathbf{a} die Projektionsrichtung, die durch Maximierung des Projektionsindex gefunden wurde, und \mathbf{U} eine orthonormale $(q \times q)$ -Matrix mit \mathbf{a} in der ersten Zeile. Durch die Lineartransformation

$$\mathbf{T} = (T_1, \dots, T_q)^\top = \mathbf{U} \mathbf{z} \quad (9.20)$$

erhält man eine Rotation, sodass die erste Koordinate

$$T_1 = \mathbf{a}^\top \mathbf{z} = y \quad (9.21)$$

ist. Sei weiters $\Theta = (\theta_1, \dots, \theta_q)^\top$ eine Transformation, deren Komponenten durch

$$\begin{aligned} \theta_1(T_1) &= \Phi^{-1}(F_a(T_1)) \ , \\ \theta_j(T_j) &= T_j \quad \text{für } 2 \leq j \leq q \end{aligned} \quad (9.22)$$

definiert sind. $F_a(T_1)$ ist die kummulative Verteilungsfunktion der eindimensionalen Projektion (9.21). Durch Anwendung einer Transformation mit der Inversen der kummulative Standardnormalverteilungsfunktion erhält man als Resultat eine Standardnormalverteilung. Die Transformation Θ wandelt also T_1 in eine Standardnormalverteilung um und lässt die Komponenten T_2, \dots, T_q unverändert.

Mit der Transformation

$$\mathbf{z}^* = \mathbf{U}^\top \Theta(\mathbf{U}\mathbf{z}) \quad (9.23)$$

ist die durch die Projektion $y = \mathbf{a}^\top \mathbf{z}$ gefundene Struktur entfernt und die Maximierung des Projektionsindex kann auf \mathbf{z}^* wieder angewandt werden.

Für konkrete Stichproben werden die theoretischen Größen von (9.22) durch die geschätzten Größen ersetzt und man erhält

$$y_i^* = \Phi^{-1}\left(\hat{F}_a(y_i)\right) = \Phi^{-1}\left(\frac{r(y_i) - \frac{1}{2}}{n}\right) \quad \text{für } i = 1, \dots, n \quad (9.24)$$

$\hat{F}_a(y_i)$ ist die empirische Verteilungsfunktion und $r(y_i)$ ist der Rang von y_i über die n projizierten Beobachtungen. Die Transformation (9.24) wird mit *Gaussianisierung* bezeichnet, sie ersetzt jede Beobachtung durch ihren normalverteilten Wert in der Projektion.

Es muss beachtet werden, dass wiederholte Gaussianisierung die Normalverteilung von nichtorthogonalen früher gefundenen Projektionen verändert, sodass der Projektionsindex größer als Null wird. Die so fälschlicherweise erhaltene Struktur ist aber in den meisten Fällen gering.

Zur Entfernung von zweidimensionalen Strukturen benötigt man eine Transformation, die eine allgemeine bivariate Dichtefunktion $p_{ab}(y_1, y_2)$ in die Dichte der bivariaten Normalverteilung

$$\phi(y_1, y_2) = \frac{1}{2\pi} e^{-\frac{y_1^2 + y_2^2}{2}} \quad (9.25)$$

umwandelt. Es ist bekannt, dass jede eindimensionale Projektion der bivariaten Normalverteilung normalverteilt ist. Daher muss $p_{ab}(y_1, y_2)$ in alle Richtungen des zweidimensionalen Unterraums auf Normalverteilung transformiert werden. Praktisch wäre das nicht durchführbar, daher normalisiert man nur in ein paar Richtungen, die man durch Rotation erhält.

Sei

$$\begin{aligned} y_1^* &= y_1 \cos \gamma + y_2 \sin \gamma \\ y_2^* &= y_2 \cos \gamma - y_1 \sin \gamma \end{aligned} \quad (9.26)$$

eine Rotation um den Winkel γ . Die Verteilungen von y_1^* und y_2^* werden auf die selbe Weise wie im eindimensionalen Fall zur Normalverteilung transformiert. Dieser Vorgang kann für verschiedene Werte von γ ($0, \frac{\pi}{8}, \frac{\pi}{4}, \frac{3\pi}{8}$) so lange wiederholt werden, bis die Verteilungen hinreichend genau die Normalverteilung approximieren.

Seien \mathbf{a} und \mathbf{b} die Projektionsrichtungen, die den Projektionsindex maximieren, und sei \mathbf{U} eine orthonormale $(q \times q)$ -Matrix mit \mathbf{a} und \mathbf{b} in den ersten beiden Zeilen. Mit der Transformation

$$\mathbf{T} = \mathbf{U}\mathbf{z} \quad (9.27)$$

erhält man als die ersten beiden Koordinaten von \mathbf{T} die Koordinaten der Projektionsebene. Mit einer weiteren Transformation $\boldsymbol{\Theta} = (\theta_1, \dots, \theta_q)^\top$, definiert durch

$$\begin{aligned} \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} \begin{pmatrix} T_1 \\ T_2 \end{pmatrix} &= \Phi^{-1}(F_{ab}(T_1, T_2)) \quad , \\ \theta_j(T_j) &= T_j \quad \text{für } 3 \leq j \leq q \quad , \end{aligned} \quad (9.28)$$

wird erreicht, dass durch

$$\mathbf{z}^* = \mathbf{U}^\top \boldsymbol{\Theta}(\mathbf{U}\mathbf{z}) \quad (9.29)$$

die Lösungsebene $\{\mathbf{a}, \mathbf{b}\}$ von PP in eine bivariate Standardnormalverteilung umgewandelt wird und alle orthogonalen Projektionen unverändert bleiben.

9.2.4 Robustheit

Der ein- und zweidimensionale Projektionsindex ist robust gegenüber extremen Ausreißern, weil das Maß für die Abweichung von der Normalverteilung hauptsächlich durch das Zentrum der Verteilung beeinflusst wird und nur gering durch die Schwänze. Multivariate Ausreißer, die keine extremen Werte in irgendeine Richtung annehmen, beeinflussen den Projektionsindex allerdings schon. Solche Ausreißer können aber möglicherweise mit PP entdeckt werden, wenn sie eine starke Struktur liefern, sodass der Projektionsindex maximiert wird. Die Ausreißer können für die weitere Analyse ausgeschlossen werden.

Die Entfernung einer Struktur ist durch Ausreißer nicht beeinflusst. Nur die Skalierung der Daten ist nicht robust, da sie auf dem Stichprobenmittel und der Stichprobenkovarianzmatrix basiert und beide stark von Ausreißern beeinflusst werden. Allerdings gibt es bereits Verfahren für robuste Skalierung wie den MVE-Schätzer (*minimum volume ellipsoid*) und den MCD-Schätzer (*minimum covariance determinant*).

9.3 Andere Projektionsindices

9.3.1 Der Projektionsindex von Friedman und Tukey

Friedman und Tukey (1974) entwickelten einen Projektionsindex, um Cluster in den Daten zu finden. Dabei wird ein Kriterium maximiert, sodass möglichst viele Punkte

in einem Cluster liegen, während die Cluster möglichst separiert sein sollten. Der Projektionsindex hat somit die Form

$$I(\mathbf{a}) = s(\mathbf{a}) d(\mathbf{a}) ,$$

\mathbf{a} ist die Projektionsachse, $s(\mathbf{a})$ misst die Streuung der Daten und $d(\mathbf{a})$ beschreibt eine „lokale Dichte“ der Punkte nach erfolgter Projektion auf \mathbf{a} . Für $s(\mathbf{a})$ wurde (aus Gründen der Robustheit) eine getrimmte Standardabweichung und für $d(\mathbf{a})$ eine „durchschnittliche Nähefunktion“ herangezogen.

9.3.2 Der Entropieindex

Sei f eine Dichtefunktion und $\int f \ln f$ ein Entropiemaß. Es kann gezeigt werden, dass jene Dichte mit Mittel Null und Varianz Eins, die dieses Maß minimiert, die Dichte der Standardnormalverteilung ist. Nimmt man also $\int f \ln f$ als Projektionsindex, so wird dadurch die Abweichung von der Dichte der Standardnormalverteilung gemessen. Diese Darstellung wurde von Huber (1985) eingeführt. Er zeigte auch, dass die *Fisher-Information* $\int \frac{(f')^2}{f}$ eine andere Möglichkeit für einen Projektionsindex ist, der durch die Dichte der Normalverteilung optimiert wird.

Wenn man einen Entropieindex auf eine konkrete Stichprobe anwenden möchte, muss vorher die Dichte f geschätzt werden. Dies kann durch eine *Kerndichteschätzung* (siehe z.B. Scott, 1992)

$$\hat{f}(y) = \frac{1}{nw} \sum_{i=1}^n K\left(\frac{y - y_i}{w}\right) \quad y \in \mathbb{R} \quad (9.30)$$

geschehen. K ist eine Kernfunktion, $w \geq 0$ ein *Smoothing*-Parameter und y_i ist der i -te Datenpunkt, der auf \mathbf{a} projiziert wurde. K ist typischerweise eine symmetrische nichtnegative Funktion mit Integral Eins, die monoton gegen Null geht, wenn ihr Argument gegen unendlich strebt. Ein Beispiel dafür ist die Dichte der Standardnormalverteilung.

9.3.3 Der Momentindex

Jones und Sibson (1987) zeigten, dass das Entropiemaß $\int f \ln f$ durch das dritte und vierte Moment, μ_3 und μ_4 , in der Form

$$\int f \ln f \approx \frac{\mu_3^2 + \frac{(\mu_4 - 3)^2}{4}}{12} \quad (9.31)$$

approximiert werden kann. Für zentrierte und skalierte normalverteilte Daten ist dieser Ausdruck Null, und er kann daher als Projektionsindex verwendet werden.

Ein Index, der auf dem dritten und vierten Moment basiert, wird möglicherweise auch durch andere Verteilungen optimiert. Allerdings sind bei konkreten Daten die höheren Momente in der Regel nicht sehr klein.

Literatur

- J.H. Friedman. Exploratory Projection Pursuit. *J. Amer. Statist. Assoc.*, 82(397):249–266, 1987.
- J.H. Friedman and W. Stuetzle. Projection Pursuit Regression. *J. Amer. Statist. Assoc.*, 76(376):817–823, 1981.
- J.H. Friedman and J.W. Tukey. A Projection Pursuit Algorithm for Exploratory Data Analysis. *IEEE Transactions on Computers*, c-23(9):881–890, 1974.
- J.H. Friedman, W. Stuetzle, and A. Schroeder. Projection Pursuit Density Estimation. *J. Amer. Statist. Assoc.*, 79(387):599–608, 1984.
- P.J. Huber. Projection Pursuit. *The Annals of Statistics*, 13(2):435–475, 1985.
- M.C. Jones and R. Sibson. What is Projection Pursuit? *J. Roy. Statist. Soc. A*, 150(1):1–36, 1987.
- D.W. Scott. *Multivariate Density Estimation: Theory, Practice, and Visualization*. John Wiley & Sons, New York, 1992.