

Multivariates Testen & Clusteranalyse

Aufgabe 1: MANOVA

In einer kalifornischen Studie wurden in 409 Schuldistrikten die mittleren Lese- und Mathefähigkeiten aller Grundschüler erhoben. Mithilfe einer multivariaten Varianzanalyse soll nun überprüft werden, wie sich die Fähigkeiten der Schüler unterscheiden, wenn man die Distrikte anhand ihres kategorisierten Einkommensniveaus vergleicht (73 Distrikte haben ein niedriges Niveau, 280 ein mittleres, 67 ein hohes).

- a) Notieren Sie die Modellgleichung des MANOVA-Modells in Matrixform. Geben Sie dabei insbesondere den Aufbau und die Dimensionen der einzelnen Matrizen an.
- b) Durch die Streuungszerlegung erhalten Sie die Matrizen

$$\mathbf{B} = \begin{pmatrix} 71446 & 63907 \\ 63907 & 57855 \end{pmatrix}, \quad \mathbf{W} = \begin{pmatrix} 97968 & 81919 \\ 81919 & 89514 \end{pmatrix}.$$

Testen Sie auf Basis von Wilks' Λ und mit einem Signifikanzniveau von $\alpha = 0.05$, ob sich die Gruppen signifikant unterscheiden.

Hinweis: Quantile von Wilks' Λ -Verteilung werden üblicherweise nur über Approximationen bestimmt, z.B. anhand von

$$\left(\frac{p - n + 1}{2} - m \right) \log \Lambda(p, m, n) \sim \chi_{np}^2.$$

- c) Welche alternativen Möglichkeiten hätten Sie, den möglichen Gruppenunterschied auf Basis einer MANOVA zu testen?
- d) Vollziehen Sie den berechneten Test in R nach:
 - (i) Laden Sie den Datensatz `CASchools` aus dem Paket `AER` und kategorisieren Sie das mittlere Einkommen in der Spalte `income` in die drei Kategorien $\{(5, 10], (10, 20], (20, 60]\}$.
 - (ii) Verschaffen Sie sich einen Überblick über die Daten (die relevanten Zielgroessen sind `read` und `math`). Überprüfen Sie insbesondere, ob die in der MANOVA gewählte Annahme der Varianzhomogenität für die vorliegenden Daten sinnvoll erscheint.
 - (iii) Schätzen Sie mithilfe der Funktion `stats::manova` die MANOVA und testen Sie mittels `stats::summary.manova` die Fragestellung unter Verwendung des Signifikanzniveaus von $\alpha = 0.05$. Wie können Sie auswählen, welcher Test durchgeführt wird?

Aufgabe 2: Hierarchische Clusteranalyse

Für vier Filialen einer Supermarktkette erhält man für die Merkmale Umsatz und Verkaufsfläche, jeweils gemessen in geeigneten Einheiten, die folgende Datenmatrix:

Filiale	1	2	3	4
Umsatz	8	5	10	4
Verkaufsfläche	24	22	25	21

- a) Führen Sie ein hierarchisches Clustering mit dem *Single Linkage* Verfahren durch. Verwenden Sie als zugrundeliegende Distanz zwischen einzelnen Objekten die quadrierte euklidische Distanz.
- b) Führen Sie ein hierarchisches Clustering mit dem *Zentroid* Verfahren durch.
- c) Geben Sie für beide Verfahren das vollständige Dendrogramm an.

Aufgabe 3: Clusteranalyse in R I

Der Datensatz `europa.txt` enthält Daten zu $n = 24$ europäischen Ländern. Folgende Variablen wurden erhoben: `ober` (Oberfläche in km^2), `einw` (Einwohner in Millionen), `brut` (BIP pro Kopf in \$) und `arbl` (Arbeitslosenquote in %).

- a) Lesen Sie den Datensatz in R ein und standardisieren Sie die Daten.
- b) Führen Sie mithilfe der Funktion `hclust()` ein hierarchisches Clustering unter Einbeziehung aller vier Kovariablen mit dem *Single Linkage* Verfahren durch. Verwenden Sie als zugrundeliegende Distanz zwischen einzelnen Objekten die quadrierte euklidische Distanz.
- c) Führen Sie mithilfe der Funktion `hclust()` ein hierarchisches Clustering unter Einbeziehung aller vier Kovariablen mit dem *Zentroid* Verfahren durch.
- d) Führen Sie mithilfe der Funktion `hclust()` ein hierarchisches Clustering unter Einbeziehung aller vier Kovariablen mit dem *Complete Linkage* Verfahren durch. Verwenden Sie als zugrundeliegende Distanz zwischen einzelnen Objekten die Mahalanobis-Distanz.

Hinweis: Die Funktion `mahalanobis()` könnte hilfreich sein.

- e) Visualisieren und vergleichen Sie ihre Ergebnisse der Teilaufgaben b), c) und d) jeweils mithilfe eines Dendrogramms.
- f) Führen Sie das k-means Clusterverfahren mit Hilfe der Funktion `kmeans()` aus dem Paket `cluster` durch. Wählen Sie dafür $k = 2, \dots, 10$ und entscheiden Sie anschließend mittels des *elbow criterion* bezogen auf die Summe der quadratischen Abstände innerhalb der Cluster was eine geeignete Anzahl an Clustern sein könnte.
- g) Wiederholen Sie die Verfahren aus den Teilaufgaben b), c), d) und f) nur unter Einbeziehung der beiden Variablen `arbl` und `brut` durch. Vergleichen Sie die Ergebnisse für $k=4$, indem Sie jeweils die 4 Cluster in 4 verschiedenen Farben im zweidimensionalen Raum plotten.