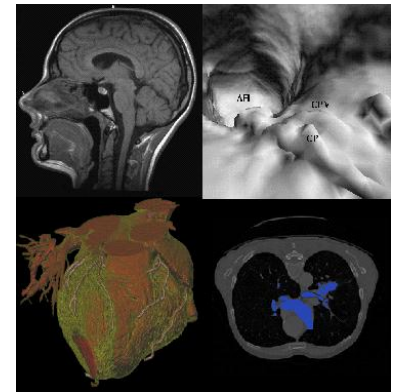
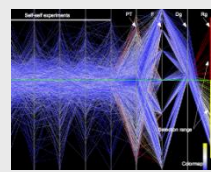


Informations- visualisierung

Thema: 3. Statistische Grundlagen
Dozent: Prof. Dr. Gerek Scheuermann
scheuermann@informatik.uni-leipzig.de
Sprechstunde: nach Vorlesung
Umfang: 2
Prüfungsfach: Modul Fortgeschrittene Computergrafik
Medizininformatik, Angewandte Informatik





Übersicht

1. Einführung
2. Wahrnehmung von Grafik
3. **Statistische Grundlagen**
4. Darstellung von Tabellen
5. Darstellung von Graphen
6. Darstellung von Metadaten und Prozessen
7. Interaktion
8. Visual Analytics



3.0 Multivariate Statistik

Informationsvisualisierung zielt auf das **Erkennen von Zusammenhängen** zwischen Variablen bei einer großen Menge von gegebenen Werten oder die **Überprüfung vermuteter Zusammenhänge**. Dazu erfolgt eine Abbildung der Variablen in graphischen Darstellungen, um das visuelle System des Menschen zur Zusammenhangssuche verwenden zu können.

Die **klassischen Ansätze** zur Analyse von mehreren Variablen bei einer großen Menge von Werten gehören in den Bereich der **multivariaten Statistik**. Da häufig gerade das Erkennen unbekannter Zusammenhänge in der Informationsvisualisierung eine Rolle spielt, wird in diesem Kapitel ein wichtigstes **strukturerkennendes Verfahren** der multivariaten Statistik vorgestellt, da es zuweilen als Ergänzung oder Vorbereitung von Informationsvisualisierungen eingesetzt wird.

Es sei darauf hingewiesen, dass hier stets von **ganzzahligen oder reellen Variablen** ausgegangen wird. Es gibt durchaus auch statistische Ansätze für ordinale oder nominale Variablen!



3.0 Multivariate Statistik

- Bevor wir zu den einzelnen Verfahren kommen, benötigen wir eine geeignete Notation, die auf Matrizen beruht.
- Wenn wir m **Datensätze (Beobachtungen)** d_1, \dots, d_m mit n **Komponenten (Variablen)** x_1, \dots, x_n haben, so bilden wir die **Datenmatrix (Beobachtungsmatrix)** X :

$$X := \begin{pmatrix} x_{11} & \dots & x_{1n} \\ \cdot & \dots & \cdot \\ \cdot & \dots & \cdot \\ x_{m1} & \dots & x_{mn} \end{pmatrix}$$

3.0 Multivariate Statistik

Mit Hilfe der Mittelwerte $\bar{x}_j := \frac{1}{m} \sum_{i=1}^m x_{ij}$

bestimmen wir die **Abweichungsmatrix** M

$$m_{ij} := x_{ij} - \bar{x}_j \quad \text{mit} \quad \begin{array}{l} i = 1, \dots, m \\ j = 1, \dots, n \end{array}$$

was sich mit Hilfe des Mittelwertvektors $\bar{\mathbf{x}}$ und des m -Vektors $\mathbf{1}$

$$\bar{\mathbf{x}} := (\bar{x}_1, \dots, \bar{x}_n) \quad \mathbf{1} := (\mathbf{1}, \dots, \mathbf{1})^T$$

als $M = X - \frac{1}{m} (\mathbf{1}\mathbf{1}^T X)$ schreiben lässt.

3.0 Multivariate Statistik

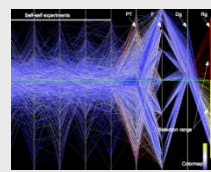
Daraus leiten wir die **Abweichungsquadratmatrix** T ab:

$$\begin{aligned} t_{jk} &:= \sum_i m_{ij} m_{ik} & j, k = 1, \dots, n \\ &= \sum_i (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \end{aligned}$$

also

$$T = M^T M = \begin{pmatrix} t_{11} & \dots & t_{1n} \\ \cdot & & \cdot \\ t_{n1} & \dots & t_{nn} \end{pmatrix}$$

die aufgrund ihrer Konstruktion symmetrisch ist.



3.0 Multivariate Statistik

Es gilt:

$$\begin{aligned} T &= M^T M = \left(X^T - \frac{1}{m} (X^T \mathbf{1} \mathbf{1}^T) \right) \left(X - \frac{1}{m} (\mathbf{1} \mathbf{1}^T X) \right) \\ &= X^T X - \frac{1}{m} X^T \mathbf{1} \mathbf{1}^T X = X^T \left(I - \frac{1}{m} \mathbf{1} \mathbf{1}^T \right) X = X^T H X \end{aligned}$$

mit

$$H := I - \frac{1}{m} \mathbf{1} \mathbf{1}^T = \begin{pmatrix} 1 - \frac{1}{m} & -\frac{1}{m} & \dots & -\frac{1}{m} \\ -\frac{1}{m} & 1 - \frac{1}{m} & \cdot & \cdot \\ \cdot & \cdot & \cdot & -\frac{1}{m} \\ -\frac{1}{m} & \dots & -\frac{1}{m} & 1 - \frac{1}{m} \end{pmatrix}$$

3.0 Multivariate Statistik

Nun kann man die **Kovarianzmatrix** S durch

$$s_{jk} := \frac{1}{m} t_{jk} = \frac{1}{m} \sum (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k) \quad j, k = 1, \dots, n$$

also

$$S := \frac{1}{m} T = \frac{1}{m} X^T H X$$

festlegen.

Einige Verfahren nutzen die **Korrelationsmatrix** R

$$r_{jk} := \frac{t_{jk}}{\sqrt{t_{jj}} \sqrt{t_{kk}}} = \frac{s_{jk}}{\sqrt{s_{jj}} \sqrt{s_{kk}}}$$

die mit Hilfe der Diagonalmatrix

$$\text{diag}T = \begin{pmatrix} t_{11} & & \\ & \cdot & \\ & & t_{nn} \end{pmatrix}$$

als $R = (\text{diag}T)^{-0.5} T (\text{diag}T)^{-0.5}$ geschrieben werden kann.

3.0 Multivariate Statistik

Häufig wird auch von der standardisierten Datenmatrix Z

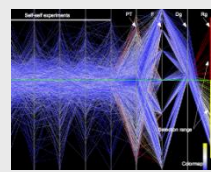
$$z_{ij} := \frac{x_{ij} - \bar{x}_j}{\sqrt{m-1}s_{jj}}$$

also

$$Z := \frac{1}{\sqrt{m-1}} M(\text{diag}T)^{0.5}$$

ausgegangen.

Es gilt $R = Z^T Z$

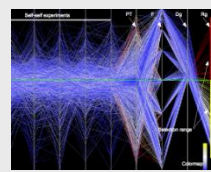


3.1 Hauptkomponentenanalyse

Mit der **Hauptkomponentenanalyse** (principle component analysis, PCA) kann man die n untereinander korrelierten beobachtbaren Variablen x_1, \dots, x_n auf möglichst wenige $p < n$ untereinander **unkorrelierte Variablen zurückführen**, die **möglichst viel Varianz** auf sich vereinigen.

Dazu werden die standardisierten Variablen z_i in neue Variablen f_i , die sogenannten **Hauptkomponenten**, linear transformiert, die **untereinander unkorreliert und nach fallender Varianz geordnet** sind. Die ersten $p < n$ Hauptkomponenten vereinigen ein Maximum der Gesamtvarianz auf sich.

[Anmerkung: In der Bildverarbeitung gibt es eine entsprechende Bildtransformation. Die Transformation heißt nach ihren Erfindern Karhunen und Loeve auch KL-Transformation.]



3.1 Hauptkomponentenanalyse

Man sucht eine Darstellung der standardisierten Datenmatrix $Z \in R^{mn}$ in der Form

$$Z = FL^T$$

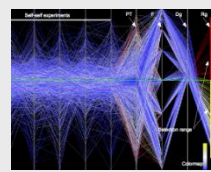
mit der **Faktormatrix** $F \in R^{mn}$ mit orthonormierten Spalten f_i und der **Ladematrix** $L \in R^{mn}$. Es gilt

$$R = Z^T Z = LF^T FL^T = LL^T \quad (F \text{ orthogonale Matrix})$$

Wählt man die $p < n$ ersten normierten Hauptachsen f_1, \dots, f_p aus, die die wichtigsten „Faktoren“ sind, so nähert man Z an:

$$Z \approx F^{(p)} L^{(p)T} \text{ mit } Z \in R^{mn}, F^{(p)} \in R^{mp} \text{ orthogonal}, L^{(p)} \in R^{np}$$

Dabei nimmt man an, dass die n Variablen x_1, \dots, x_n linear unabhängig sind, was durch Weglassen von Spalten in X bzw. Z stets erreicht werden kann.



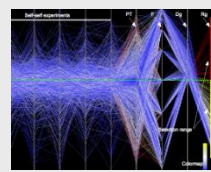
3.1 Hauptkomponentenanalyse

Im Wesentlichen erfordert dies eine **Hauptachsentransformation** der Korrelationsmatrix R . Da R symmetrisch und reell ist, gibt es n Eigenwerte $\lambda_1, \dots, \lambda_n$ und orthogonale Eigenvektoren e_1, \dots, e_n mit

$$Re_j = \lambda_j e_j, \quad |e_j| = 1$$

also

$$R = E\Lambda E^T \quad \text{mit} \quad E = (e_1 \dots e_n), \quad \Lambda := \begin{pmatrix} \lambda_1 & & \\ & \cdot & \\ & & \lambda_n \end{pmatrix}$$



3.1 Hauptkomponentenanalyse

Die Matrix der Hauptachsen Y wird definiert durch

$$Y := Z E$$

und es gilt

$$Y^T Y = E^T Z^T Z E = E^T R E = \Lambda$$

Setzt man nun

$$F := Y \Lambda^{-0.5}, \quad L := E \Lambda^{-0.5}$$

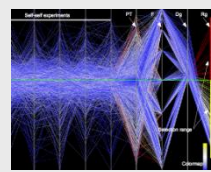
so hat man

$$F^T F = I, \quad Z = F L^T, \quad R = Z^T Z = L L^T$$

Dabei heißt die Zerlegung

$$Z = F^T \Lambda^{0.5} E^T$$

auch **Singulärwertzerlegung** von Z .



3.1 Hauptkomponentenanalyse

Geometrische Deutung

Seien z_1, z_2 gemeinsam normalverteilte Zufallsvariablen, also

$$\begin{pmatrix} z_1 \\ z_2 \end{pmatrix} \approx N(\mathbf{0}, \Sigma), \quad \Sigma = \begin{pmatrix} 1 & \sigma \\ \sigma & 1 \end{pmatrix}$$

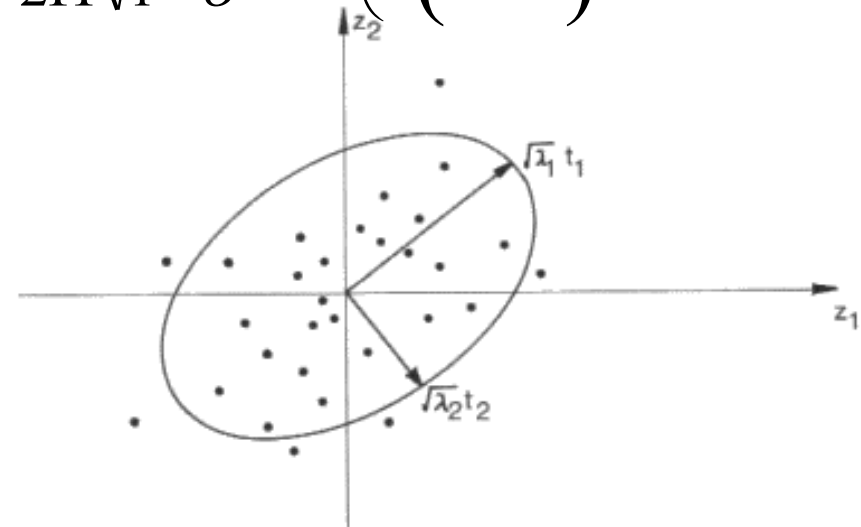
Dann liegen Punkte gleicher Dichte

$$f(z_1, z_2) = \frac{1}{2\pi\sqrt{1-\sigma^2}} \exp\left(\frac{-1}{2(1-\sigma^2)} (z_1^2 - 2\sigma z_1 z_2 + z_2^2)\right)$$

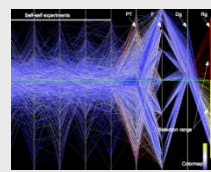
auf Ellipsen um den Nullpunkt.

M Beobachtungen (hier $M=30$)

geben die elliptische Form wieder:



Die Hauptachse t_1 gibt die Richtung größter Streuung, t_2 die Reststreuung an.



3.1 Hauptkomponentenanalyse

Komplette Analyse

Zunächst muss entschieden werden, ob überhaupt Abhängigkeiten in den Variablen vorliegen. Dazu stellt man die Hypothese H_0 auf:

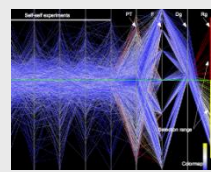
$$H_0: R = I \quad H_1: R \neq I$$

Unter der Annahme, dass es sich bei Z um normalverteilte Zufallsvariablen handelt, kann man dann einen χ^2 -Test durchführen.

Dazu vergleicht man

$$\chi_{err1}^2 := -(m - (2n + 11)/6) \ln(\det R)$$

mit den Werten der χ^2 -Verteilung. Nach Festlegung der Irrtumswahrscheinlichkeit α wird also ermittelt, ob sich R von I signifikant unterscheidet.



3.1 Hauptkomponentenanalyse

Wenn R von I verschieden ist, erfolgt die Hauptachsentransformation und die Darstellung

$$\mathbf{Z} = \mathbf{F}\mathbf{L}^T$$

Um noch die Anzahl $p < n$ der wichtigen Hauptkomponenten zu bestimmen, notiert man zunächst, dass der Anteil der Hauptkomponente f_1 an der Gesamtvarianz n gleich λ_1/n ist.

Mit Hilfe der Testmaßzahl
$$\chi_{err2}^2 := (m-1)(n-p) \ln \det R$$

prüft man die Hypothese

H_2 : Die $n-p$ kleinsten Eigenwerte von R sind gleich

Als Faustregel können stattdessen einfach die Hauptkomponenten mit $\lambda_j > 1$ benutzt werden. Beide Methoden sind nicht perfekt, da auch im ersten Fall χ_{err2}^2 für $n \rightarrow \infty$ nicht gegen χ^2 konvergiert.