

Lehrstuhl fuer Statistik, Oekonometrie und empirische  
Wirtschaftsforschung, Universitaet Tuebingen

# Hauptkomponentenanalyse

Dr. S. Prohl

8. Juni 2007/ SoSe 2007

# Hauptkomponentenanalyse

- ▶ Problemstellung der Hauptkomponentenanalyse
- ▶ Hauptkomponentenanalyse bei bekannter Varianz-Kovarianz-Matrix
- ▶ Hauptkomponentenanalyse bei unbekannter Varianz-Kovarianz-Matrix
  - ▶ Anzahl der Hauptkomponenten
  - ▶ Überprüfung der Güte der Anpassung
- ▶ Beispiel

# Hauptkomponentenanalyse

## Literatur:

- ▶ Handl, A.: Multivariate Verfahren: Theorie und Praxis multivariater Verfahren unter besonderer Beruecksichtigung von S-Plus, 2002, Kapitel 5.

# Hauptkomponentenanalyse: Vorgehensweise

- ▶ Ausgangslage: hochdimensionaler Datensatz, der in einem niedrig-dimensionalen Raum dargestellt werden soll.
  - ▶ Alle Merkmale sind quantitativ.
  - ▶ Die Daten liegen als Varianz-Kovarianz-Matrix oder Korrelationsmatrix vor.
- ▶ Soll die Hauptkomponentenanalyse auf Basis der Varianz-Kovarianz-Matrix, oder auf Basis der Korrelationsmatrix durchgeführt werden?
- ▶ Man bestimmt die Eigenwerte und Eigenvektoren der Varianz-Kovarianz-Matrix bzw. der Korrelationsmatrix.
- ▶ Wie viele Hauptkomponenten benötigt man?
- ▶ Die Hauptkomponenten werden interpretiert.

# Hauptkomponentenanalyse: Ausgangslage

Ausgangsgroessen:

Zufallsvektor

$$\mathbf{x}^T = (x_1, \dots, x_p)$$

$$\boldsymbol{\mu} = E(\mathbf{x})$$

$$\boldsymbol{\Sigma} = \text{Cov}(\mathbf{x})$$

bzw. Datenmatrix  $\mathbf{x}_1, \dots, \mathbf{x}_n$  mit  $\bar{\mathbf{x}}$  und empirischer Kovarianzmatrix  $\boldsymbol{\Sigma}$ .

# Hauptkomponentenanalyse: Problemstellung

- ▶ Reduziere die Dimension von  $\mathbf{x}$  mit moeglichst geringem Informationsverlust
- ▶  $\Rightarrow$  Es wird gesucht:
  - ▶ Linear-Kombination  $y = \mathbf{a}^T \mathbf{x}$  mit grosstem  $V(y) = \mathbf{a}^T \Sigma \mathbf{a}$

# Hauptkomponentenanalyse: Loesung des Problems

Erste Hauptkomponente

- ▶ Finde die Loesung:

$$a^T \Sigma a \quad \text{unter der Nebenbedingung} \quad a^T a = 1$$

- ▶ Die Lagrange-Funktion lautet:

$$\begin{aligned} L(a, \lambda) &= a^T \Sigma a - \lambda (a^T a - 1) \\ \frac{\partial L(a, \lambda)}{\partial a} &= 2 \Sigma a - 2 \lambda a \\ \frac{\partial L(a, \lambda)}{\partial \lambda} &= 1 - a^T a \\ \Rightarrow (\Sigma - \lambda) a &= 0 \end{aligned}$$

- ▶  $\lambda$  ist Eigenwert zu  $\Sigma$ ,  $a$  ist Eigenvektor

# Hauptkomponentenanalyse: Loesung des Problems

Es folgt

$$a^T \Sigma a = a \lambda a = \lambda$$

Es handelt sich um ein Eigenwertproblem!

- ▶ Da wir eine Linearkombination mit der groessten Varianz suchen, waehlen wir den Eigenvektor, der zum groessten Eigenwert  $\lambda_1$  gehoert.
- ▶ Erste Hauptkomponente ist  $y_1 = a_1^T x$ , und  $a_1$  ist der Eigenvektor zum groessten Eigenwert von  $\Sigma$ .

# Hauptkomponentenanalyse: Loesung des Problems

Zweite Hauptkomponente:

- ▶ Finde die Loesung:

$\mathbf{a}_2^T \Sigma \mathbf{a}_2$  unter der Nebenbedingungen  $\mathbf{a}_1^T \mathbf{a}_2 = 0$

- ▶ Die Lagrange-Funktion lautet:

$$L(\mathbf{a}_2, \lambda, \nu) = \mathbf{a}_2^T \Sigma \mathbf{a}_2 - \lambda (\mathbf{a}_2^T \mathbf{a}_2 - 1) - \nu (\mathbf{a}_2^T \mathbf{a}_1)$$

$$\frac{\partial L(\mathbf{a}_2, \lambda, \nu)}{\partial \mathbf{a}_2} = 2(\Sigma - \lambda) \mathbf{a}_2 - \nu \mathbf{a}_1 = 0.$$

$$\Rightarrow (\Sigma - \lambda) \mathbf{a}_2 = 0.$$

- ▶  $\mathbf{a}_2$  ist orthogonal zu  $\mathbf{a}_1$ .
- ▶  $\mathbf{a}_2$  ist zweite Hauptkomponente.

# Hauptkomponentenanalyse: Zahl der Hauptkomponenten

Wie viele Hauptkomponenten benoetigt man?

- ▶ Anteil der Gesamtstreuung, die durch Hauptkomponenten erklart wird,
- ▶ Kaiser-Kriterium ,
- ▶ Jolliffe-Kriterium,
- ▶ Scree-Plot