

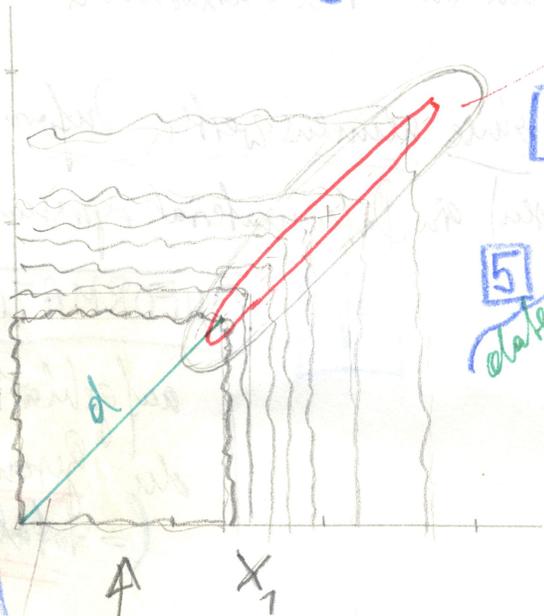
Hauptkomponentenanalyse, PCA

(principal components analysis)

(A) Zweck: Datenreduktion > einfachere Darstellung d. ursprünglich ermittelte Daten ohne Messwertverlust Daten d. ursprünglichen Informationen

Variablen X_1, X_2

1 gegeben ist eine Stichprobe zum Beispiel von



3 Centroid (\bar{x}_1 / \bar{x}_2)

falls es sich um eine sehr flache Punkt Wolke handelt, > ausstelle einer Vermessung von x_1 und x_2 könnte man, um eine Größe die die Größe d. Vermessen:

$$d = \sqrt{x_1^2 + x_2^2}$$

5 datenreduzierender Effekt

(quadratische Briefmarken als Beispiel ähnlich) Fliesen, Dachziegel, Schildkrötenpanzerplatten

6

stellt eine **LINEAR KOMBINATION** dar

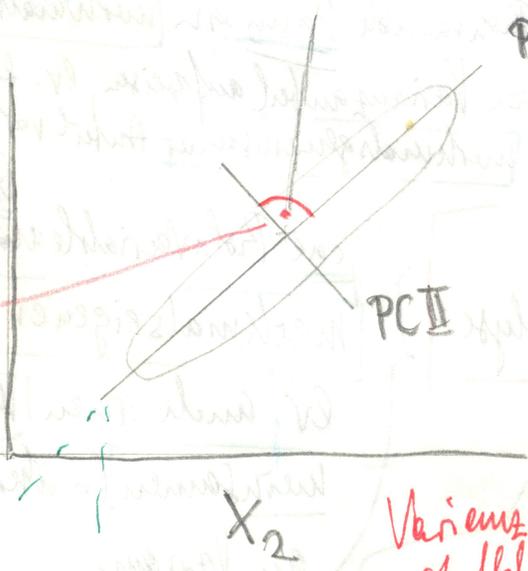
(linear compound, linear combination)

(B) Vorgehensweise der PCA: Centroid (\bar{x}_1 / \bar{x}_2)

PC... principal Component = Hauptachse

ORTHOGONALITÄT X_2

Die PC I und PC II stehen aufeinander normal



PC I **7**

Es wird eine neues Koordinatensystem nach folgenden Kriterien gesucht:

(1) die 1. Achse soll ein Maximum an Varianz d. Punkt wolke aufweisen [Varianz d. PC I ist Eigenwert-1]

(2) die folgende(n) Achse erklart von der verbleibenden Varianz wieder ein Maximum ...

(3) die neuen Achsen stehen

Varianz-staffelung

Varianz d. Komponente = Eigenwert (Richtung wird Eigenvektor angegeben)

aufeinander normal (\Rightarrow **ORTHOGONAL**), d.h. die neuen Achse = neue Variablen, sind **nicht** miteinander **korreliert**. $\hat{=}$ **UNKORRELIERTHEIT**

8

(4) Wenn die Rohdaten miteinander korreliert waren, führt mit der PCA ein datenreduzierender Effekt zu erzielen, d.h. anstatt, dass Sie ursprünglich z.B. 10 Rohvariablen zur Datenrepräsentation benötigen können ev. 2 Hauptachsen (PCs) ausreichen, um die Daten darzustellen

der datenred. Effekt liegt darin, dass ohne numerischen Informationsverlust mit geringeren Aufwand (Achsen) ein Datenmaterial repräsentiert werden kann

9

Ausgang der PCA von

10

S oder R
(Var.-Cov.-Matrix) (Korrelationsmatrix)

9

INFORMATION VOM:
aufgeklart wird die Gesamtvarianz (= Total Variance)

11

✓ der unterschiedliche Ausgang bedingt, dass die Ergebnisse d. PCA in der Regel unterschiedlich sind

(1) Die PCA ist auch Grundlage für **Faktorenanalyse** (factor analysis)

Konzeptionelle Differenzen Zwischen

PCA =

{ eine Rohvariable kann einen merkmalseigenen und einen merkmalgemeinsamen Anteil aufweisen, ev. sind aber nur ein merkmalgemeinsamer Anteil vorhanden

(orthogonale) Faktorenanalyse:

jede Rohvariable besitzt einen merkmalseigenen Anteil und ev. auch einen merkmalgemeinsamen (= überlappenden) Anteil an Varianz

merkmalseig. Anteil: E_i -- der Variable i (Error)

$\hat{=}$ exklusiv, unpartialisiert $\hat{=}$ unkorreliert mit allen anderen "Errors"

(E) Beispiel aus der Morphometrie:

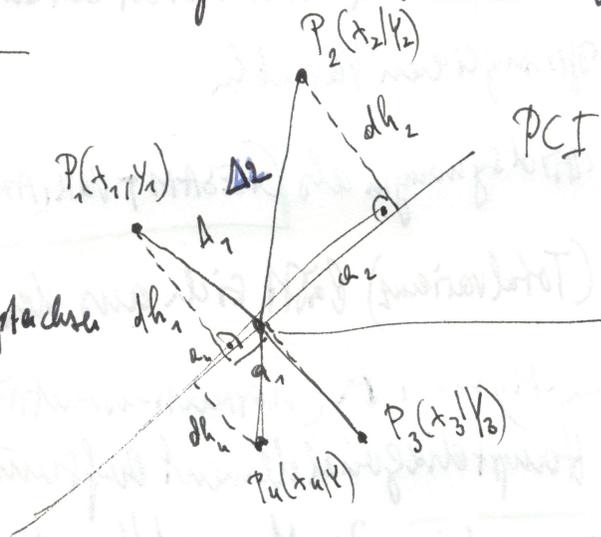
z.B. >40 Skelettmaße, transexuell erhoben,
logarithmisch transformiert, lassen sich auf 2 Hauptachsen
reduzieren > Aufklärungsgrad $\approx 98\%$ der Totalvarianz

- 1. Achse: sog. "Größenachse" ('size axis')
- 2. Achse: sog. "Formachse" ('shape axis')

Nachtrag ad (B): Nachtrag > orthogonale Regression als alternativer
Ansatz für die Hauptkomponentenanalyse

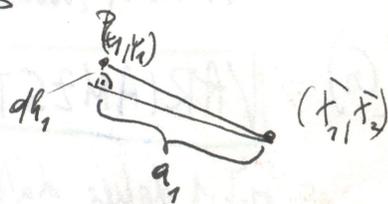
Koordinaten:
① ursprüngl. Koordinate $P(x_i | y_i)$

② Koordinaten auf d. Hauptachsen $P(a_i | dh_i)$



$(\bar{x}_1, \bar{x}_2, \dots)$ Centroid

minimales dh_1 :



a) es gilt d. pythagor. Lehrsatz in d. Ebene:

$$\Delta_1^2 = dh_1^2 + a_1^2$$

b) wenn gilt: $\sum_{i=1}^n dh_i^2 \rightarrow \min$

dann gilt: $dh_i^2 = \Delta_i^2 - a_i^2$

da Δ_i^2 als Centroidentfernung eines Punktes
konstant ist > wird a_i^2 maximiert

\rightarrow folglich die Varianz auf PC1
maximal!

(F) a) Die Varianz d. Hauptachsen werden **EIGENWERTE** (eigenvalues) genannt symbol $\lambda \dots$ lambda d.h. $\lambda_1, \lambda_2, \dots, \lambda_m$

b) die Richtungen der Hauptachsen werden durch die **EIGENVEKTOREN** (eigenvectors) angegeben

c) Die Summe aller Eigenwerte eines Systems entspricht der **TOTALVARIANZ** (total variance) d. Systems

d) Die TOTALVARIANZ entspricht zum anderen der Summe der Varianz d. ursprünglichen Variablen

e) Die Totalvarianz wird synonym als **GESAMTVARIANZ** bezeichnet.

f) Die Gesamtvarianz (Totalvarianz) lässt sich aus dem Matrizen S (Varianz-Kovarianz-M.) und R (Korrelationsmatrix) leicht ermitteln, indem man die Hauptdiagonalelemente aufsummiert

(G) **VARIANZSTAFFELUNG**: Der Name 'Hauptkomponenten-Analyse'

führt daher, das im Zuge d. Verfahrensverlaufes beginnenden mit den wichtigsten Achsen (d.h. Achsen, welche die größte ursprüngliche Varianz darstellen) gefolgt von den weniger wichtigen Achsen (mit geringerer Varianzaufklärung)

die Achsenextraktion durchgeführt wird bis bei einer vollständigen Extraktion die am wenigsten wichtige Achse am Ende kommt.

Der Extraktionsfortschritt ist gestaffelt.

Eine VOLLSTÄNDIGE EXTRAKTION ergibt VARIANZSTAFFELUNG,

die Anzahl der neuen Achsen (Hauptachsen) stimmt aber Zahlenmäßig mit der Anzahl d. ursprünglichen Achsen (= ROHACHSEN) überein.

Extraktionsfolge: $PC I \rightarrow PC II \rightarrow PC III \dots PC_m$ bei m ROHACHSEN

VARIANZSTAFFELUNG: $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \lambda_m$