

PCA 1

Zur Erläuterung und theoretischen Vertiefung der Methode
der Hauptkomponentenanalyse (= PCA) sollte die Durchdringung
viele konkreten Beispielen hilfreich sein.

In der Folge werden 2 Beispiele präsentiert, zuerst das 1.

Beispiel 1: umfangreiches nach Anzahl der Merkmalsvariablen
z. die nach Zahl der zugrundeliegenden Arten und Individuen ist.



Beispiel-1 'BOHNEN' für PCA:

- a) Das Datenmaterial stammt von einer Stichprobenziehung
einer Bohnensorte und hatte den Umfang $n = 50$
[siehe CBo.xls]

Aus Gründen allometrischer Zusammenhänge war zu vermuten,
dass die zu messenden Merkmalsvariablen ohne Transformation
nicht primär linear zueinander bezogen sind.

Aus diesem Grund wurden die abgenommenen Messdaten d. Variablen logarithmiert (zur Basis e, log. naturalis).

Als Messvariablen wurden ursprünglich LÄNGE, BREITE und GEWICHT abgenommen, mit den physikalischen Dimensionen LÄNGE [mm], BREITE [mm] und GEWICHT [g]. Nach der logarithmischen Transformation liegen die Variablen als LLG, LBR und LGEW vor.

Die logarithmische Transformation führt im Fall allometrischer Zusammenhänge zu einer LINEARISIERUNG der Datenzusammenhänge, welche GRUNDVORAUSSETZUNG für das ALM (= Allgemeine Lineare MODELL) sind.

- b) Als Grundlage für die nachfolgende Durchführung einer PCA wird zuerst eine Produkt-Moment-Korrelationsmatrix R zwischen d. 3 Merkmalsvariablen berechnet [siehe KORRMA.PDF]. Der die Rohvariablen in unterschiedlichen physikalischen Messdimensionen, u. z.B. $[m \cdot m]$ und $[g]$ erhoben werden, das für die Auszähl der zugrundeliegenden Matrix d. PCA nur die Korrelationsmatrix aufgrund der damit erzielten Dimensionalität d. Parameter möglich] WICHTIG! Der Produkt-Moment-Korrelationskoeffizient führt zur Normierung d. Variablen und zu deren Dimensionalität bezogen auf die physikalische Dimension)

- c) Das Protokoll für die VOLLSTÄNDIGE LÖSUNG der PCA findet sich unter RLBOPCA.pdf.

Nachstehende Punkte sind zu beachten:

- ① Der Ausgang nimmt das Verfahren von der Korrelationsmatrix (3-dimensional, KORRMA.PDF), $n=50$ (individuen als OTVs bezeichnet).
- ② Die Gesamtvarianz (= 'total variance') kann im Falle einer zugrundeliegenden Korrelationsmatrix leicht angegeben werden: $tv = \text{total variance} ; t_v = \text{Anzahl d. Merkmalsvariablen} = \text{Summe d. Hauptdiagonale in der Korrelationsmatrix}$
- ③ Auf die Gesamtvarianz bezieht sich d. Gesamte Variancebilanzierung d. PCA.
- ④ Das Auffinden d. 1. Hauptachse (Hauptkomponente, PC-1) bedeutet, den Hauptanteil d. Gesamtvarianz

insofern ausfindig zu machen, als eine Achse zu finden ist, welche diesen Hauptanteil repräsentiert.

Man kann auch formulieren: Es ist die Achse im Merkmalsraum zu finden, welche ein MAXIMUM an der GESAMTVARIANZ darstellt.

- ⑤ Die Varians d. neuen Achse im Merkmalsraum wird als EIGENWERT (= 'eigenvalue') bezeichnet.
- ⑥ Die Position d. neuen Achse im Merkmalsraum wird ('eigen vector') durch den EIGENVEKTOR bestimmt.
- ⑦ Die sukzessive extraktierten Hauptachsen erklären in der Folge immer geringere Anteile an d. Gesamtvarianz
[VARIANZSTAFFELUNG $\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots$]
- ⑧ Wenn man stattdessen vom wormierten Eigenvektor ausgingen, den mit der Standardabweichung der Komponenten $\sqrt{\lambda_i}$ skalierten Eigenvektor zugrunde legt, so erhält man einen sog. FAKTOREN(EIGEN)-VEKTOR. Die KOEFFIZIENTEN dieses Vektors werden als LADUNGEN bezeichnet.
- ⑨ Die Loadungen sind einer sehr einfachen Interpretation zugänglich. Es handelt sich um EINFACHE KORRELATIONSKOEFFIZIENTEN der jeweiligen Rohvariablen (= Variable, die d. Analyse zugrunde liegen) mit d. Hauptkomponente.

10

Das Quadrat dieser FAKTOREN LADUNGEN (= LDG QUA) entspricht d. erfachten BESTIMMTHEITSMASS und gibt damit Auskunft darüber, welche Anteile der Variationen d. einzelnen Rohvariablen in einer Komponente repräsentiert sind.
($\hat{=}$ AMV AR)

11

Es lässt sich ebenso bilanzieren, welchen Anteil eine bestimmte Rohvariable an der Aufklärung der Gesamtvarianz an einer bestimmten Hauptkomponente hat [$\hat{=}$ MAG VAR]

d) Die Interpretation d. vorliegenden Beispiels soll folgende Punkte umbedingt umfassen:

①

Die 1. Hauptachse wird durch einen Eigenvektor charakterisiert, dessen Koeffizienten alle am Werte mit positivem Vorzeichen und sehr ähnlicher Mächtigkeit zeigen. Variieren auf dieser PC-1 bedeutet, dass die sich zusammensetzenden Rohvariablen auf dieser Achse GLEICH SINNIG variieren. Das bedeutet, wenn z.B. die Länge einer Bohne zunimmt, so handelt es sich um Bohnen, deren Breite ebenso wie deren Gewicht zunehmen. Umgeshoben, wenn z.B. die Länge einer Bohne zunimmt, f.d. ist zu erwarten, dass es sich um eine Bohne handelt, deren Breite und deren Gewicht ebenso abnehmen. Die PC-1 ist als sog. 'Großen-Variations-Achse' ($\hat{=}$ size axis) zu interpretieren.

(2)

Die 1. Hauptachse und ihre Interpretationen haben aufgrund der Tatsache, dass diese 1. Hauptachse den „Löwenanteil“ an der Gesamtvarianz, nämlich ~ 78%, repräsentiert, das größte Gewicht (im Sinne von Bedeutung für die Interpretation) Hingegen haben alle folgenden Achsen weit weniger Gewicht.

(3)

Die PC-2 wird durch einen Eigenvektor charakterisiert, welcher durch eine Vorzeichenkontrastierung der Koeffizienten für Länge und Breite auffällt. Die niedrige absolute Koeffizientenwert für das Gezeit erlaubt es diesen Koeffizienten außer Acht zu lassen, welchen an dieser Achse offenbar keine wesentliche Bedeutung zukommt.

Auf dieser Achse variiieren die Bohnen der Stichprobe gegengleich, auf LÄNGE und BREITE, bezogen:
D.h., findet man Bohnen im Sample mit größerer Länge, so ziehen sich diese Bohnen auch durch geringere Breite aus (> 'schlanke Bohnen').

Umgekehrt gilt: Bohnen mit geringerer Länge sind im Sample eher breiter (> 'dickere Bohnen').

(4)

Interessant ist die Koeffizientencharakteristik des 3. Eigenvektors. Hier werden LÄNGE und BREITE zum GEWICHT kontrastiert. D.h. Bei (wenigen) Bohnen der Stichprobe findet man Individuen, die sehr groß erscheinen (LÄNGE u. BREITE), aber geringeres Gezeit aufweisen. Umgekehrt gibt es auch kleinere Bohnen, die relativ schwer sind.

e) Als letzter Teil d. Analysen im RLBOPCA.pdf findet sich unter 'Ordnation der Daten aus der faktorenmétrischen Matrix' die Representation der Koordinatenwerte für die Daten der 50 Individuen auf den neuen Achsen (= 3 Hauptachsen).