

## Kapitel 4 Dimensionsreduktion

### § 1 Hauptkomponentenanalyse (PCA)

PCA  $\hat{=}$  principal component analysis

Gegeben: Datenpunkte  $X = \left( \begin{array}{c} x_1^1 & x_1^2 & \dots & x_1^n \\ \vdots & \vdots & & \vdots \end{array} \right) \in \mathbb{R}^{D \times N}$

Oft:  $D$  sehr groß.

Hoffnung: alle  $x^n$  liegen nahe zu einer niedrig-dimensionalen Mannigfaltigkeit, d. h. die Daten  $x^n$  haben Struktur.

einfaches Modell: alle  $x^n$  liegen nahe zu einem affinen Unterraum der Dimension  $M$  und  $M \ll D$ .

D.h.

$$x^n \approx c + \sum_{i=1}^M y_i^n b_i^i = \tilde{x}_n, \text{ wobei } b_1, \dots, b_M \in \mathbb{R}^D \text{ linear unabhängig.}$$

Ziel: Finde  $c \in \mathbb{R}^D$ ,  $Y = (y_i^n)_{i,n} \in \mathbb{R}^{M \times N}$ ,  $B = [b_1^T \dots b_M^T]^T \in \mathbb{R}^{D \times M}$ ,

so dass

$$E[B, Y, c] = \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2.$$

minimal ist

Vereinfachung: Datenpunkte sind zentriert, d.h. es gilt

$$\frac{1}{N} \sum_{n=1}^N x^n = 0.$$

In diesem Fall kann man zeigen, dass  $c=0$  sein muss  
(→ Blatt 9).

Satz Ang. die Datenpunkte  $(\bar{x}^n)$  sind zentriert.

Seien  $\lambda_1, \dots, \lambda_D$  die Eigenwerte der Matrix  $X X^T \in S_{\geq 0}^D$ ,  
so sortiert, dass  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$  gilt. Sei  
 $b_1, \dots, b_D$  eine Orthonormalbasis von  $\mathbb{R}^D$  bestehend aus  
zugehörigen Eigenvektoren. Dann gilt

$$\min_{B, Y} E[B, Y] = \sum_{i=M+1}^D \lambda_i$$

und ein optimales  $B$  ist  $B = [b_1 \ | \ b_M]$  und  $Y = B^T X$ .  
(Die Vektoren  $b_1, \dots, b_M$  heißen Hauptkomponenten des  
Datensatz  $X$ .)

Bew.: Es gilt

$$E[B, Y] = \sum_{n=1}^N \|x^n - \tilde{x}^n\|^2 = \sum_{n=1}^N \sum_{i=1}^D (x_i^n - \tilde{x}_i^n)^2 = \dots$$

$$\begin{aligned}
 &= \sum_{n=1}^N \sum_{i=1}^D \left( x_i^n - \sum_{j=1}^M y_j^n b_i^j \right)^2 \\
 &= \text{Tr} \left( (X - BY)^T (X - BY) \right) \quad (*) 
 \end{aligned}$$

Wir können oBdA annehmen, dass  $B$  aus orthonormalen Vektoren besteht, d.h. dass  $B^T B = I_M$  gilt.

Für fixes  $B$  haben wir folgendes notwendige Kriterium für die Optimalität von  $Y$ :

$$\begin{aligned}
 0 &= \frac{\partial}{\partial y_{jk}^n} E[B, Y] = \frac{\partial}{\partial y_{jk}^n} \sum_{i=1}^D \left( x_i^n - \sum_{j=1}^M y_j^n b_i^j \right)^2 \\
 &= \sum_{i=1}^D \frac{\partial}{\partial y_{jk}^n} \left( (x_i^n)^2 - 2x_i^n \sum_{j=1}^M y_j^n b_i^j + \sum_{j,j'=1}^M y_j^n y_{j'}^n b_i^j b_i^{j'} \right) \\
 &= \sum_{i=1}^D (-2x_i^n b_i^k) + 2 \sum_{j=1}^M y_j^n b_i^j b_i^k \\
 &= -2 \sum_{i=1}^D x_i^n b_i^k + 2 \sum_{j=1}^M y_j^n \underbrace{\sum_{i=1}^D b_i^j b_i^k}_{\delta_{jk}} \\
 &= -2 \sum_{i=1}^D x_i^n b_i^k + 2 y_k^n \\
 \Rightarrow & \boxed{Y = BX} \quad (\star\star)
 \end{aligned}$$

Setze nun  $(\star\star)$  in  $(\star)$  ein:

$$\begin{aligned} E[B] &= \text{Tr}((X - BB^T X)^T (X - BB^T X)) \\ &= \text{Tr}((X^T - X^T BB^T)(X - BB^T X)) \\ &= \text{Tr}(X^T X - X^T BB^T X - X^T BB^T X \\ &\quad + \underbrace{X^T BB^T BB^T X}_{I_n}) \\ &= \text{Tr}(X^T X - X^T BB^T X) \\ &= \text{Tr}(X^T(I - BB^T)X) \\ &= \text{Tr}((I - BB^T)XX^T) \\ &= \text{Tr}(XX^T) - \text{Tr}(BB^T XX^T) \rightarrow \min. \end{aligned}$$

DR.  $\text{Tr}(BB^T XX^T) \rightarrow \max.$ , unter der Nebenbedingung  
 $B^T B = I$ .

Aus der VL Konvexe Optimierung (Lecture 10)

Theorem (Fan)

Sei  $C \in S^n$  sym. Matrix mit Eigenwerten  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ .

Dann gilt

$$\lambda_1 + \dots + \lambda_k = \max \left\{ \text{Tr}(C BB^T) : B \in \mathbb{R}^{n \times k}, B^T B = I_k \right\}$$

$$= \max \left\{ \text{Tr}(CY) : Y \in S^n, \text{Tr}(Y) = k, Y \succeq 0, I - Y \succeq 0 \right\}.$$

Wende dies auf  $C = XX^T$  an. Dann

$$\begin{aligned} E[B] &= \text{Tr}(XX^T) - \text{Tr}(BB^TXX^T) \\ &= \sum_{i=1}^D \lambda_i - \sum_{i=1}^M \lambda_i = \sum_{i=M+1}^D \lambda_i. \end{aligned}$$

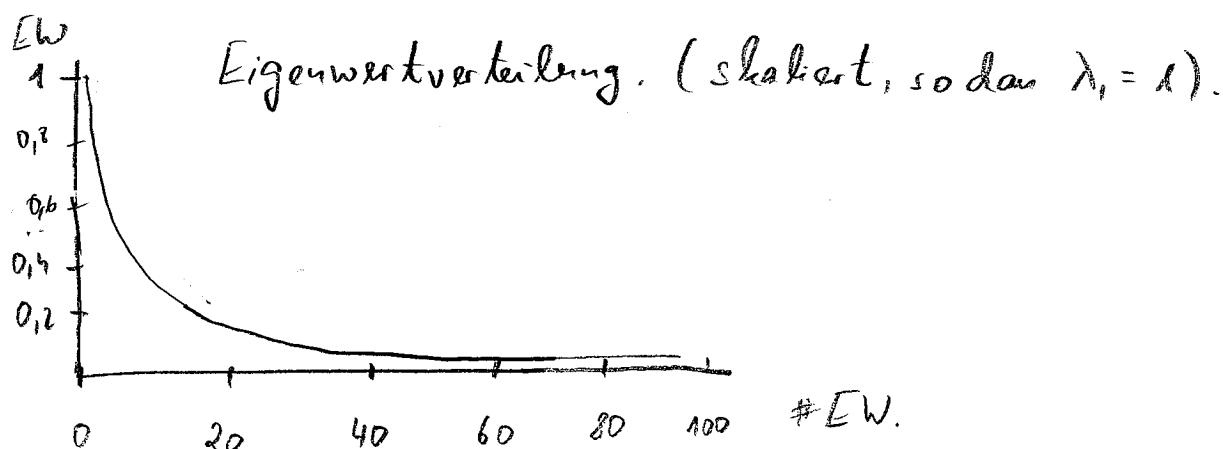
Dass  $B$  aus den zu  $\lambda_1, \dots, \lambda_D$  gehörigen Eigenvektoren von  $XX^T$  besteht, folgt ebenfalls aus dem Beweis des Theorems von Fan.  $\square$

Bsp.: (für die Hauptkomponentenanalyse)

$X \in \mathbb{R}^{784 \times 892}$  Datenmatrix der handgeschriebenen 5  
der MNIST-Datenbank

$N = 892$  Anzahl Beispiele

$D = 784 = 28 \times 28$  Pixel pro Bruchstabe.



## PCA: Beispiel aus MNIST Datenbank (aus Barber)

