

## Vorkurs 19 PCA cont'd

Erinnerung: Wir haben Daten  $z_1, \dots, z_n \in \mathbb{R}^M$  und  $d < M$ .

In PCA suchen wir einen linearen Raum  $U \subset \mathbb{R}^M$  von Dim.  $d$  und  $b \in \mathbb{R}^M$ , s.d.

$$U + b \subset \mathbb{R}^M$$

"nah" an den Daten.

Unser erstes Konzept von "nah" war ein Raum von maximaler Varianz.

Dazu haben wir angenommen, dass  $z_1, \dots, z_n$  unabh. Samples einer ZV  $z \in \mathbb{R}^M$ . Dann war

$$U = \text{span}\{u_1, \dots, u_d\}, \quad b = \mu,$$

wobei  $\mu = \mathbb{E}z$ ,  $u_1, \dots, u_d$  Eigenvektoren der Kovarianzmatrix  $\Sigma$  von  $z$  zu den größten Eigenwerten.

In der Praxis approximieren wir  $\Sigma$  durch die empirische Kovarianzmatrix  $S$  und  $\mu$  durch den empirischen Mittelwert  $\bar{z}$ .

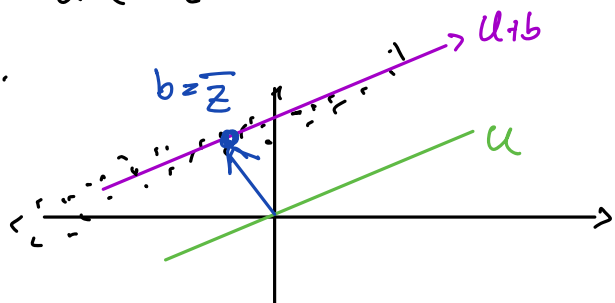
Bemerkung:  $u_1, \dots, u_d$  heißen Hauptkomponenten ("Principal Components")

Nun betrachten wir ein zweites Konzept von "nah".

Wir wählen  $U$ , s.d.

$$\alpha) \quad \sum_{i=1}^n \|(z_i - \bar{z}) - P_U(z_i - \bar{z})\|^2$$

minimiert wird und  $b = \bar{z}$ .



### Theorem 18.1

Seien  $\lambda_1, \dots, \lambda_M \geq 0$  die Eigenwerte der emp. Kovarianzmatrix  $S$ .

Sei  $u_i$  Eigenvektor zu  $\lambda_i$ , sd.  $\langle u_i, u_j \rangle = \delta_{ij}$ . Dann

minimiert

$$U = \text{span}\{u_1, \dots, u_d\}$$

den Ausdruck in (\*). Falls  $\lambda_d > \lambda_{d+1}$ , ist  $U$  eindeutig bestimmt.

Beweis Sei  $u_1, \dots, u_M$  eine Orthonormal-Basis von  $\mathbb{R}^M$ , d.h.  $\langle u_i, u_j \rangle = \delta_{ij}$ .

Sei  $w_i = z_i - \bar{z}$ . Sei außerdem  $A = [u_1, \dots, u_d] \in \mathbb{R}^{M \times d}$

und  $W = [w_1 \dots w_n] \in \mathbb{R}^{M \times n}$

Wir haben außerdem die Feature-Matrix  $\Omega = [z_1, \dots, z_n]^T \in \mathbb{R}^{n \times M}$ .

Dann gilt:

$$W = \Omega^T - \bar{z} e^T, \quad e = (1, \dots, 1) \in \mathbb{R}^n.$$

Dann gilt:

$$\begin{aligned} WW^T &= (\Omega^T - \bar{z} e^T)(\Omega^T - \bar{z} e^T)^T \\ &= nS. \end{aligned}$$

Außerdem:  $AA^T = \sum_{i=1}^d u_i u_i^T$  und  $I_M = \sum_{i=1}^M u_i u_i^T$

Dann gilt:

$$\begin{aligned} W - P_U W &= W - AA^T W \\ &= (I_M - AA^T) W \\ &= \left( \sum_{i=d+1}^M u_i u_i^T \right) W \end{aligned}$$

Außerdem,  $W - P_U W = [w_1 - P_U(w_1), \dots, w_n - P_U(w_n)]$ .

Es gilt:

$$\begin{aligned}
(*) &= \text{Trace} \left( (W - P_U W)^T (W - P_U W) \right) \\
&= \text{Trace} \left( W^T \left( \sum_{i=d+1}^M u_i u_i^T \right) \left( \sum_{i=d+1}^M u_i u_i^T \right) W \right) \\
&= \text{Trace} \left( W^T \left( \sum_{i=d+1}^M u_i u_i^T \right) W \right) \quad , \text{ weil } \langle u_i, u_j \rangle = \delta_{ij} \\
&= \sum_{i=d+1}^M \text{Trace} (W^T u_i u_i^T W) \\
&= \sum_{i=d+1}^M \text{Trace} (u_i^T W W^T u_i) \\
&= n \cdot \sum_{i=d+1}^M u_i^T S u_i
\end{aligned}$$

Wie im Beweis von Thm. 18.4 erhalten wir:

$u_1, \dots, u_d$  sind Eigenvektoren zu  $\lambda_1, \dots, \lambda_d$ .

Eindeutigkeit folgt ebenfalls wie in Thm. 18.4.  $\square$ .

In beiden Ansätzen müssen wir eine Eigenzerlegung von  $S = \frac{1}{n} W W^T$  berechnen. Es ist  $S \in \mathbb{R}^{M \times M}$ .  $W \in \mathbb{R}^{M \times n}$

Falls  $n \ll M$ , können wir wie folgt vorgehen. Wir nehmen an, dass  $r(W) = n$ . Dann sei

$$W = U D V^T$$

ein SVD von  $W$ , d.h.  $U \in \mathbb{R}^{M \times n}$ ,  $D = \text{diag}(\sigma_1, \dots, \sigma_n)$ ,  $V \in \mathbb{R}^{n \times n}$ .

Dann gilt:

$$S = \frac{1}{n} W W^T = U D^2 U^T$$

und die Diagonaleinträge von  $D^2$  = nicht-null Eigenwerte von  $S$ .

Außerdem:  $\mathbb{R}^{n \times n} \ni W^T W = V D^2 V^T$

$\Rightarrow$  eine Eigenzerlegung von  $\frac{1}{n}W^T W$  gibt uns die gleichen Eigenwerte wie für  $S$ .

Außerdem  $u_i = \frac{W v_i}{\|W v_i\|}$ , so dass wir auch die Hauptkomponenten von  $S$  erhalten

Insbesondere zeigt dies, dass PCA durch SVD verstanden werden kann.

Außerdem kann  $W^T W$  durch die Kernel-Map  $K(x_i, x_j) = \langle z_i, z_j \rangle = \langle \phi(x_i), \phi(x_j) \rangle$  berechnet werden.

### Lemma 18.2

Sei  $G = (K(x_i, x_j))_{i,j=1}^n$  für Daten  $x_1, \dots, x_n \in \mathbb{R}^D$ ,  $z_i = \phi(x_i)$ .

Dann gilt für  $W = [z_1 - \bar{z}, \dots, z_n - \bar{z}]$  wie oben, dass

$$W^T W = \left( I_n - \frac{1}{n} e e^T \right) G \left( I_n - \frac{1}{n} e e^T \right).$$

### Beweis

Es ist  $W = \Omega Z^T - \bar{z} e^T$  und  $\bar{z} = \frac{1}{n} (z_1 + \dots + z_n) = \frac{1}{n} \Omega Z^T e$

$$\begin{aligned} \Rightarrow W^T W &= \left( \Omega Z^T - \frac{1}{n} \Omega Z^T e e^T \right)^T \left( \Omega Z^T - \frac{1}{n} \Omega Z^T e e^T \right) \\ &= \left( \Omega Z^T \left( I_n - \frac{1}{n} e e^T \right) \right)^T \left( \Omega Z^T \left( I_n - \frac{1}{n} e e^T \right) \right) \\ &= \left( I_n - \frac{1}{n} e e^T \right) \Omega \Omega^T \left( I_n - \frac{1}{n} e e^T \right) \end{aligned}$$

Es ist:  $\Omega = \begin{bmatrix} -z_1^T \\ \vdots \\ -z_n^T \end{bmatrix}$ . D.h.  $\Omega \Omega^T = (\langle z_i, z_j \rangle)_{i,j=1}^n = G$

$$\Rightarrow W^T W = \left( I_n - \frac{1}{n} e e^T \right) G \left( I_n - \frac{1}{n} e e^T \right).$$

□.

Als nächstes betrachten wir folgendes statistisches Setting für PCA.  
Wir nehmen an, dass die Daten  $x_1, \dots, x_n \in \mathbb{R}^D$  von folgender Zufallsvariable gesampled sind:

$$X = A\zeta + b + \epsilon,$$

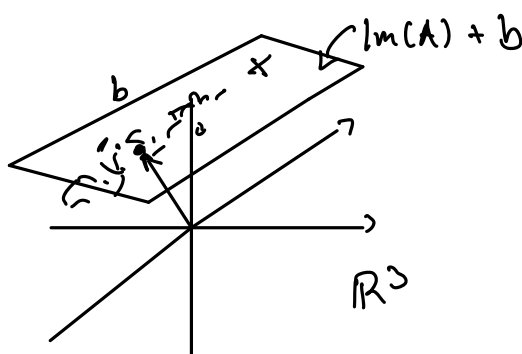
mit  $A \in \mathbb{R}^{D \times d}$ ,  $\zeta \sim N(0, 1_d)$ ,  $b \in \mathbb{R}^D$ ,  $\epsilon \sim N(0, \sigma^2 1_D)$ .  
 $\text{rank}(A) = d$ .

Im Folgenden arbeiten wir ohne Feature map.

D.h.

$$(x|\zeta) \sim N(A\zeta + b, \sigma^2 1_D).$$

Cartoon:  $d=2$ ,  $D=3$



### Proposition 18.3

Im obigen Setting haben wir

$$x \sim N(b, AA^T + \sigma^2 1_D).$$

### Beweis

Es gilt,  $P(x) = \int_{\mathbb{R}^d} P(x|\zeta) P(\zeta) d\zeta$

Es ist:

$$\log P(x|\zeta) + \log P(\zeta)$$

$$= -\frac{1}{2\sigma^2} \|x - (A\zeta + b)\|^2 - \frac{1}{2} \|\zeta\|^2 + c$$

$$= -\frac{1}{2} (x - \mu) \Sigma^{-1} (x - \mu) - \frac{1}{2} (\zeta - \gamma) B^{-1} (\zeta - \gamma) + c'$$

mit  $c, c'$  unabh. von  $x$  und  $\xi$  und  $\mu \in \mathbb{R}^D, \gamma \in \mathbb{R}^d$   
 $\Sigma \in \mathbb{R}^{D \times D}, B \in \mathbb{R}^{d \times d}$  und  $\mu$  und  $\Sigma$  unabh. von  $\xi$ .

D.h., nachdem wir  $\xi$  ausgewirkt haben, bleibt ein Ausdruck der Form  $e^{-\frac{1}{2}(x-\mu)\Sigma^{-1}(x-\mu)} \frac{1}{\sqrt{(2\pi)^D \det(\Sigma)}}$ .

$$\leadsto x \sim N(\mu, \Sigma).$$

$$\text{Es ist } \mu = \mathbb{E}_x x = \mathbb{E}_{\xi, \varepsilon} (A\xi + b + \varepsilon) = b$$

$$\Sigma = \mathbb{E}_x (x-b)(x-b)^T$$

$$= \mathbb{E}_{\xi, \varepsilon} (A\xi + b + \varepsilon - b) (A\xi + b + \varepsilon - b)^T$$

$$= \mathbb{E}_{\xi, \varepsilon} (A\xi + \varepsilon) (A\xi + \varepsilon)^T$$

$$= \mathbb{E}_{\xi, \varepsilon} (A\xi \xi^T A^T + A\xi \varepsilon^T + \varepsilon \xi^T A^T + \varepsilon \varepsilon^T)$$

$$= \mathbb{E}_{\xi, \varepsilon} (A\xi \xi^T A^T + \varepsilon \varepsilon^T).$$

$$= A \mathbb{E} \xi \xi^T A^T + \mathbb{E} \varepsilon \varepsilon^T$$

$$= A A^T + \sigma^2 \mathbb{1}_D.$$

□.

### Korollar 18.4

Angenommen wir haben den Prior  $\xi \sim N(\gamma, B)$ . Dann erhalten wir

$$x \sim N(A\gamma + b, ABA^T + \sigma^2 \mathbb{1}_D).$$

### Beweis

Es ist  $\xi = R\xi' + \gamma$  mit  $\xi' \sim N(0, \mathbb{1}_D)$ , und  $RA^T = B$ .

Daher:

$$x = A\xi + b = AR\xi' + A\gamma + b.$$

Wir wenden Thm. 18.3 mit  $AR$  und  $A\gamma + b$  an. □.

Mit Hilfe von Kor. 18.4 können wir jetzt MLE für die Parameter  $A, b$  bestimmen. Wir können auch die Verteilung von  $(\eta | x)$  ausrechnen, um den Prior von  $\eta$  nach Sicht von  $x$  upzudaten.

### Theorem 18.5

Sei:  $\eta \sim N(\nu, B)$  und  $(x | \eta) \sim N(A\eta + b, \sigma^2 I_D)$ .

Dann gilt:  $(\eta | x) \sim N(m, C)$ .

mit  $C = (\sigma^{-2} A^T A + B^{-1})^{-1}$  und  $m = C(\sigma^{-2} A^T (x - b) + B^{-1} \nu)$

### Beweis

Es gilt, nach Bayes' Theorem:

$$\log P(\eta | x) = \log P(x | \eta) + \log P(\eta) + c$$

mit  $c$  unabh. von  $\eta$  und  $x$ .

Es ist  $(x | \eta) \sim N(A\eta + b, \sigma^2 I_D)$  und  $\eta \sim N(\nu, B)$ , s.d.

$$\log P(\eta | x) = -\frac{1}{2\sigma^2} \|x - (A\eta + b)\|^2 - \frac{1}{2}(\eta - \nu)^T B^{-1}(\eta - \nu) + c''$$

mit  $c''$  unabh. von  $x$  und  $\eta$ .

$$= -\frac{1}{2\sigma^2} \|A\eta - (x - b)\|^2 - \frac{1}{2}(\eta - \nu)^T B^{-1}(\eta - \nu).$$

Das wollen wir nun nach  $\eta$  auflösen. Wir gehen wie in Thm 14.2 vor und erhalten die Formeln für  $m$  und  $C$ .  $\square$