



"dětoreč předmaty" v 4. roč.  $\left\{ \begin{array}{l} RDD \\ \cdot \text{Analyza z hlediska ... (Herman)} \\ \hookrightarrow \text{prebora' niektore' z metod z RDD} \\ \hookrightarrow \text{odporučam zapsat} \end{array} \right.$

RDD  $\left\{ \begin{array}{l} a) \text{ kombinovanie starých premenných, aby sme vytvorili nové: PCA, fakt. analyza ... - extrakcia premenných} \\ \rightarrow \text{①: lin. metódy, ②: nelineárne met. (feature extraction)} \\ b) \text{ vyber iba časti premenných - selekcia prem. (feature selection)} \end{array} \right.$

$\hookrightarrow$  Lasso, Huberová regresia, ...

$\rightarrow$  ③

- ale môže byť aj príliš vysoká  $n$  - najmä numericky

$\rightarrow$  konštrukcia podvzorky (subsampling)

$\rightarrow$  ④  
ak stáleme

: leveraging, ...

① Lineárne metódy extrakcie premenných - lineárne v dátach

Metóda hlavných komponentov (Principal component analysis - PCA)

$x_i \in \mathbb{R}^2$ , chceme znížiť dimenziu

- ale nechceme stratit prírodnú informáciu

- PCA: informáciu meria rozptylnosťou

teoretický / populačný pohľad:

$X = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$  - náh. vektor

- kovariančná matica:

$$\Sigma = \text{Cov}(X) = E[(X - E(X))(X - E(X))^T]$$

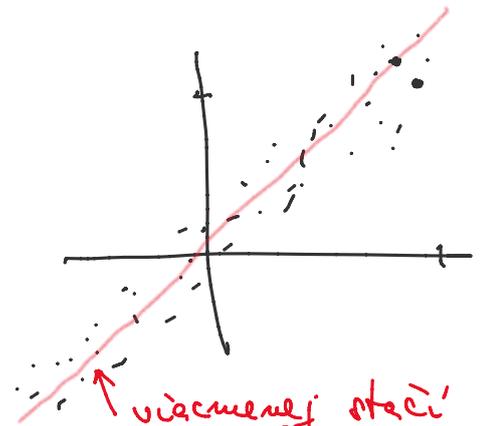
- PCA pracuje so  $\Sigma \rightarrow$  s jej spektrálnym rozkladom

$$\Sigma_{n \times n} - \text{r.s.d.} \Rightarrow \Sigma = U \Lambda U^T$$

$n \times n$       $n \times n$       $n \times n$       $n \times n$

$$U^T U = I = U U^T$$

$U$  - ortogonálna, vl. vektory = stĺpce  $U$



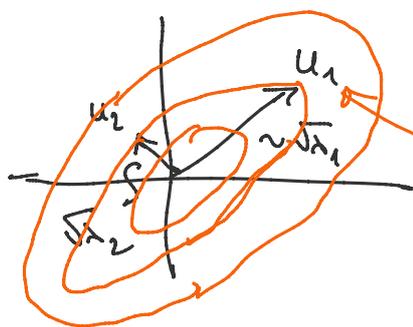
9K2

$U$  - ortogonálna, vl. vektory = stĺpce  $U$   
 $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ ,  $\lambda_j$  - vl. hodnoty  
 $\sum u_j = \lambda_j u_j \quad \forall j=1, \dots, n \quad (Ax = \lambda x)$

Príklad:  $Ax \sim N_2(0, \Sigma)$ , potom vnútornice jeho hustoty  
 sú elipsoidy so stredom v  $O_2$  a polosami v smeroch  
 vl. vektorov matice  $\Sigma$ ; dĺžky polosí  $\sim \sqrt{\lambda_j}$

$$f(x) = \text{mno.} \cdot e^{-\frac{x^T \Sigma^{-1} x}{2}} = c$$

$$\Leftrightarrow x^T \Sigma^{-1} x = \text{konst.}$$



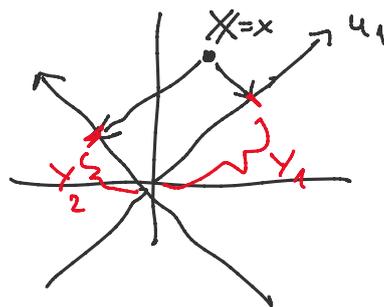
v tomto  
 smere je  
 $f(x)$   
 najviac  
 koncentrovaná

Predp.:  $X: E(X) = \mu$   
 $\text{Cov}(X) = \Sigma = U \Lambda U^T$   
 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$

- vezme sme súrad. systém daný vl. vektormi  $u_1, \dots, u_n$   
 $(U = (u_1, \dots, u_n))$
- $\rightarrow$  to sú kľúčové smery dát
- $\rightarrow$  aké sú súradnice  $X$  v súr. systéme danom  $U$ ?

$$X = U \cdot Y \quad / \quad U^T \left( X = Y_1 \cdot u_1 + Y_2 \cdot u_2 \right) \quad \text{pre } n=2$$

$$\Rightarrow U^T X = Y$$



- toto je pre  $\mu = 0$   
 pre všeob.  $\mu$ : pracovíme s  $X - \mu$ :

$$X - \mu = U \cdot Y$$

$$\Rightarrow Y = U^T (X - \mu) : \text{súradnice } X \text{ vzhľadom na } U$$

$Y_j$ : dĺžka projekcie  $X$  na priamku danú  $u_j$

$Y_j$ :  $j$ -ty hlavný komponent

Definujeme:  $Y = U^T (X - \mu)$  : vektor hlavných komponentov

$$E(Y) = U^T [E(X) - \mu] = 0_n$$

$$\text{Cov}(Y) = \text{Cov}(U^T (X - \mu)) = U^T \text{Cov}(X) U =$$

$$\boxed{\text{Cov}(AX) = A \text{Cov}(X) A^T}$$

$$\text{Cov}(Y) = \text{Cov}(U^T(X-\mu)) = U^T \text{Cov}(X) U = U^T U \Lambda U^T U = \Lambda = \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_k \end{bmatrix}_{k \times k}$$

$\Rightarrow Y_j$  sú nekorelované!

•  $\text{Cov}(X, Y) = \text{Cov}(X, U^T(X-\mu)) = \text{Cov}(X, U^T X) = \text{Cov}(X, X) \cdot U = U \Lambda U^T U = U \Lambda =$

$$= \begin{bmatrix} \lambda_1 u_{11} & \lambda_2 u_{21} & \dots & \lambda_k u_{k1} \\ \vdots & \vdots & & \vdots \\ \lambda_1 u_{1j} & \lambda_2 u_{2j} & \dots & \lambda_k u_{kj} \\ \vdots & \vdots & & \vdots \\ \lambda_1 u_{1i} & \lambda_2 u_{2i} & \dots & \lambda_k u_{ki} \\ \vdots & \vdots & & \vdots \\ \lambda_1 u_{1n} & \lambda_2 u_{2n} & \dots & \lambda_k u_{kn} \end{bmatrix}$$

(i) j-tý prvok  $U$  i-tý prvok vektora  $u_j$

$\Rightarrow \text{cov}(X_i, Y_j) = (\lambda_j u_j)_i = \lambda_j (u_j)_i = \lambda_j (u_j)_i$

$\rho_{X_i, Y_j} = \frac{\lambda_j (u_j)_i}{\sqrt{\sum_{ii} \lambda_j}} = \frac{(u_j)_i \sqrt{\lambda_j}}{\sum_{ii}}$   $\rightarrow$  j-tý vl. vektor  $u_j$  určuje korelácie medzi  $Y_j$  a  $X_1, \dots, X_n$ .

$\uparrow$   
j-tý hl. komp.

Veta: Prvý hl. komponent má najväčšiu disperziu spomedzi všetkých normovaných lin. kombinácií zložiek  $X$ . T.j.  $D(Y_1) \geq D(h^T X) \quad \forall h \in \mathbb{R}^n: \|h\|=1$ .

Dôkaz:  $D(h^T X) = h^T \text{Cov}(X) h = h^T U \Lambda U^T h = \underbrace{h^T U}_{v^T} \underbrace{\Lambda}_{\Lambda} \underbrace{U^T h}_v = v^T \Lambda v = \lambda_1 v_1^2 + \dots + \lambda_k v_k^2 \leq \lambda_1 (v_1^2 + \dots + v_k^2) = \lambda_1 \|v\|^2 = \lambda_1 v^T v = \lambda_1 h^T U U^T h = \lambda_1 h^T h = \lambda_1 = D(Y_1)$ .

$\square$

• druhý hl. komponent má najväčšiu disperziu spomedzi všetkých norm. lin. komb. na ortogonálnom komplemente k  $Y_1$ :  $D(Y_2) \geq D(h^T X) \quad \forall h \in \mathbb{R}^n: \|h\|=1, h \perp u_1$

**OPRAVA:**  $\leftarrow Y_2 = u_2^T (X - \mu)$

...komplemente ...

$$D(Y_2) \geq D(l^T X) \quad \forall l \in \mathbb{R}^2: \|l\|=1, l \perp u_1$$

OPRAVK:

$(Y_1 = u_1^T(X - \mu))$   
cize  $Y_1$  ...  
pomocou  $l = u_1$   
 $\Rightarrow$  pre dalsie  $Y_j$   
musi byt'  
 $l \perp u_1$

atd, pre  $Y_3$ :  $\rightarrow$  -

$$l \perp u_1, l \perp u_2$$