

Minule:

$$E(X) = \mu$$

PCA - teoreticky: $Cov(X) = \Sigma = U \Lambda U^T$: spektr. roz.

$Y = U^T(X - \mu)$ - vektor hl. komponentov : siradnice X vzhľadom na u_1, \dots, u_q
 $\rightarrow Y_j$ - j -ty hl. komponent

Veta: $D(Y_1) \geq D(Q^T X) \forall Q \in \mathbb{R}^q, \|Q\|=1$.

Pre $j > 1$: $D(Y_2) \geq D(Q^T X) \forall Q \in \mathbb{R}^q, \|Q\|=1, Q \perp u_1$
 $D(Y_3) \geq \dots, Q \perp u_1, Q \perp u_2$

predpokladáme: $\mu = 0$ (pre jednoduchost)

prakticky: X_1, \dots, X_n : nhl. vyber z rozdelenia s $E = \mu$ a $Var = \Sigma$
 realitacie: $x_1, \dots, x_n \in \mathbb{R}^q$

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\Sigma} = S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$
 : vyberova kovariancna matica

q_j : vl. vektory $S \rightarrow Q = (q_1, \dots, q_q)$: $q \times q$
 : ortog.

v_j : vl. hodnoty S

• ak $X_i \sim N(\mu, \Sigma)$, Σ -neg., potom $q_j \xrightarrow{n \rightarrow \infty} u_j$
 $v_j \rightarrow \lambda_j$

teoret.: $Y = U^T(X - \mu)$

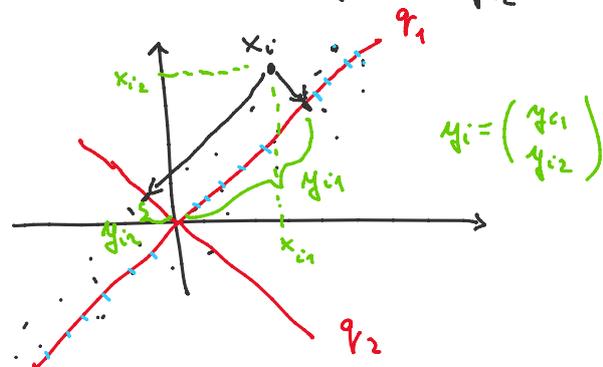
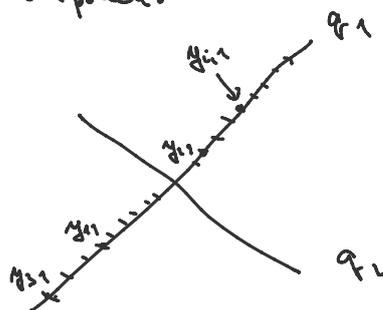
vyberovo: $x_i - \bar{x} = Q \cdot y_i$ / Q^T

$$y_i = Q^T (x_i - \bar{x}) \quad i=1, \dots, n$$

y_i - siradnice centrovaneho x_i vzhľadom na q_1, \dots, q_q

y_{11}, \dots, y_{m1} - score

$\begin{pmatrix} y_{1j} \\ y_{2j} \\ \vdots \\ y_{mj} \end{pmatrix}$ - vektor clone na j -ty hl. komponent



$$y_i = Q^T(x_i - \bar{x})$$

Předp.: x_1, \dots, x_n sú centrovane: $\bar{x} = 0_n$

↳ časť prvej časti analýzy

↳ $\tilde{x}_i = x_i - \bar{x} \quad \forall i$

↳ ocistenie prímer

$$\Rightarrow y_i = Q^T x_i \quad \forall i$$

$$(y_1, \dots, y_n) = (Q^T x_1, \dots, Q^T x_n)$$

$$X_{n \times p} = \begin{pmatrix} -x_1^T - \\ \vdots \\ -x_n^T - \end{pmatrix}$$

$$\underbrace{(y_1, \dots, y_n)}_{\text{ozn. } Y^T} = Q^T \underbrace{(x_1, \dots, x_n)}_{X^T}$$

ozn. Y^T

$$\hookrightarrow \text{ozn. } Y = \begin{pmatrix} -y_1^T - \\ \vdots \\ -y_n^T - \end{pmatrix}_{n \times p}$$

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

$$\boxed{Y = X Q}$$

$$x_1, \dots, x_n \in \mathbb{R}^p \rightarrow y_1, \dots, y_n \in \mathbb{R}^k$$

Ciel: redukovať dimenziu - vzameť iba niektoré l. komp. : tie s najvyššou variáciou

$$D(Y_1) \geq D(Y_2) \geq \dots \geq D(Y_k) \quad Y = \begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix}$$

" λ_1 " λ_2 " λ_k

→ vzameť prvých q súradníc y_1, \dots, y_q

$$\alpha_j = \frac{\lambda_1 + \dots + \lambda_j}{\lambda_1 + \dots + \lambda_k} \in [0, 1] \quad j = 1, \dots, k$$

↑ časť celkovej variácie, kt. vysvetľuje prvých j hlavných komponentov
- vysvetlená časť variácie

výberovo $\hat{\alpha}_j = \frac{v_1 + \dots + v_j}{v_1 + \dots + v_k}$

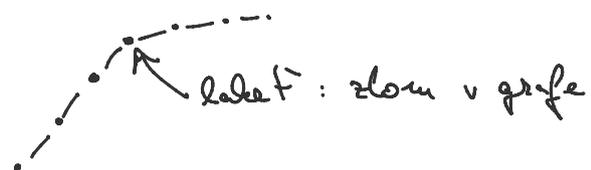
($\hat{\alpha}_j$ konverguje k α_j pre $N(\cdot)$ dát)

voľba # l. komponentov:

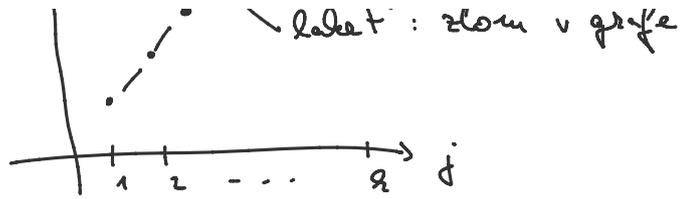
1) elbow (laktinový) diagram

- vzameť tie l. komponenty,

kt. sú do bodu zlomu.



- vezmeme tie dl. komponenty, kt. sú do bodu zlomu



2.) najmenšie j také, že $\lambda_j \geq c$
(pre zvolené c , napr. $0,9$)

3.) Kaiserovo pravidlo: najmenšie j , že $\lambda_j \geq \text{priemer}(\lambda_1, \dots, \lambda_q)$

Prk0:

food.txt - zloženie jedál

- 7. stĺpec: koľko tukov (celkový...) obsahuje na 1g jedla
- ↳ tieto uvažujeme

• PCA nie je invariantná na zmenu mierky
(ak zmeníme x na cx , zmení sa výsledok)

→ pokus o riešenie: dáta normalizujeme / štandardizujeme

- predeliť každú premennú odvozcinnou z jej disperzie ($X_i \rightarrow X_i / \sqrt{D(X_i)}$)
- potom majú všetky disperzie 1 ($D(X_i) = 1 \ \forall X_i$)
- vlastne pracujeme s $R = \cos(\alpha)$ namiesto S

→ ak sú všetky premenné rovnakého typu a v rovnakej mierke, tak to problém nie je

• mali sme: $u_j = \begin{pmatrix} u_{j1} \\ \vdots \\ u_{jq} \end{pmatrix} \rightarrow \sim \mathcal{S}(X_1, Y_j)$
 $ \phantom{\begin{pmatrix} u_{j1} \\ \vdots \\ u_{jq} \end{pmatrix}} \rightarrow \sim \mathcal{S}(X_q, Y_j)$

→ PC1: "silne" zap. korelácia s Energiou, Tuhmi a Nasyt. tuhmi

⇒ jedlo s vysokou hodnotou PC₁ (s vysokým 1. složkou)

bude menej energetické a menej mastné

→ PC2: viac uhlíkovodíkov, menej cholest. } jedlo s vysokou PC₂
viac energie, menej bielkovín } sú sčôr uhlíkovodíkové

Aplikácie PCA:

- vyčistenie, menejrozmerná reprezentácia

- vstup do iných metód: napr. zníženie # vysvetľujúcich prem. v regresii

- Principal component regression