

Unsupervised learning “učenie bez učiteľa” I. redukcia dimenzie 1.časť

Zuzana Rošťáková

23. september 2021

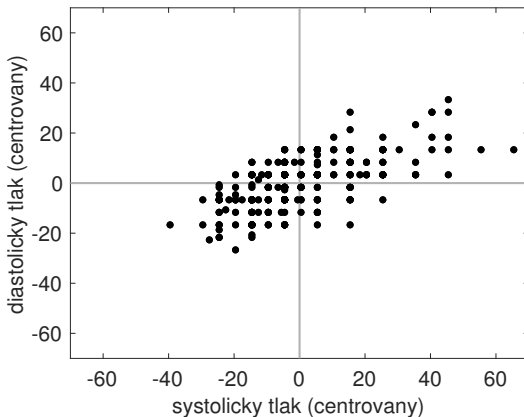
Seminár UM

- vstupné údaje: na skupine n objektov pozorujeme p premenných

$$X = \begin{bmatrix} - & \mathbf{x}_1^T & - \\ - & \mathbf{x}_2^T & - \\ & \vdots & \\ - & \mathbf{x}_n^T & - \end{bmatrix} = \begin{bmatrix} | & | & & | \\ \mathbf{x}_1^* & \mathbf{x}_2^* & \dots & \mathbf{x}_p^* \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times p}$$

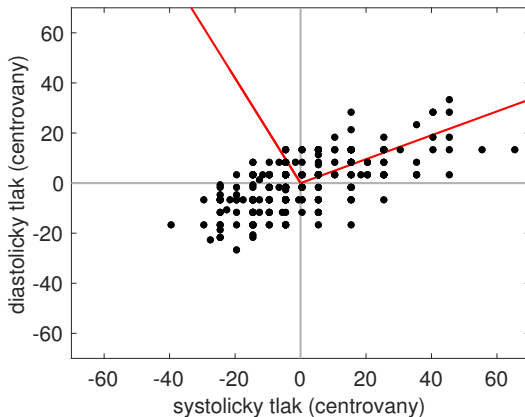
- $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$
→ vektor hodnôt p premenných pre i -ty objekt, $i = 1, \dots, n$
- $\mathbf{x}_l^* = (x_{1l}, x_{2l}, \dots, x_{nl})^T$
→ vektor hodnôt l -tej premennej, $l = 1, \dots, p$ pre n objektov

Redukcia dimenzie



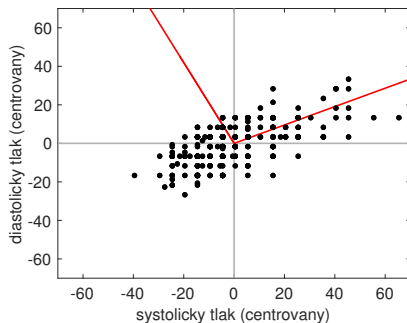
- premenné môžu byť nositeľmi “podobnej” informácie
- ⇒ stačí jedna (nová) premenná, napr. ich vhodná lineárna kombinácia

Redukcia dimenzie



- premenné môžu byť nositeľmi “podobnej” informácie
- ⇒ stačí jedna (nová) premenná, napr. ich vhodná lineárna kombinácia

Redukcia dimenzie



- **cieľ:** hľadanie K nových premenných

- $K \leq p$
- nové premenné dobre popisujú pôvodné dáta
- nové premenné vieme interpretovať

Redukcia dimenzie

$$\mathbf{X} = \mathbf{1}_n + \mathbf{S} + \mathbf{E}$$

$$x_{ij} = u_j + \sum_{k=1}^K s_{ik} c_{kj} + e_{ij}$$

- $\mathbf{S} \in \mathbb{R}^{n \times K}$ → hodnoty K “nových” premenných pre n pozorovaní
- $\mathbf{C} \in \mathbb{R}^{p \times K}$ → “vzťah” medzi K novými a p pôvodnými premennými
- $\mathbf{u} \in \mathbb{R}^p$ → vektor priemerov jednotlivých premenných

$$u_j = \frac{1}{n} \sum_{l=1}^n x_{lj}, \quad j = 1, \dots, p$$

- $\mathbf{1}_n = (1, 1, \dots, 1)^T \in \mathbb{R}^n$
- $\mathbf{E} \in \mathbb{R}^{n \times p}$ → matica chýb modelu

- odlišnosti v predpokladoch vlastností a odhadoch matíc S, C
- **metóda hlavných komponentov** (principal component analysis)
- klasická, jadrová (kernel), riedka (sparse), nelineárna, robustná ...
 - **faktorová analýza** (exploratory factor analysis)
 - **metóda nezávislých komponentov**
(independent component analysis)
- klasická, jadrová (kernel), ...
 - **nezáporná maticová faktorizácia** (nonnegative matrix factorisation)
- klasická, semi-ortogonálna, hladká, riedka, ...
 - ...

Metóda hlavných komponentov

Principal component analysis (PCA)

- **cieľ:** nájsť nové premenné, ktoré:
 - sú navzájom ortogonálne
 - sú lineárnou kombináciou pôvodných premenných
⇒ dochádza len k transformácií údajov
 - “dobre” charakterizujú dáta
⇒ vieme ich interpretovať

$$X = \mathbf{1}_n \mathbf{u}^T + SC^T + E$$

$$x_{ij} = u_j + \sum_{l=1}^p s_{il} x_{lj} + e_{ij}$$

→ $C = (\mathbf{c}_1, \dots, \mathbf{c}_p) \in \mathbb{R}^{p \times p}$

- prepojenie medzi hl. komponentami (PC, stĺpce) a pôvodnými premennými (riadky)

$$C^T C = I_p$$

→ $S = (\mathbf{s}_1, \dots, \mathbf{s}_p) \in \mathbb{R}^{n \times p}$

- "nové" premenné = hlavné komponenty
- lin. kombinácia pôvodných (centrovaných) premenných

$$s_{il} = \sum_{j=1}^p c_{lj} (x_{ij} - u_j) \quad l = 1, \dots, p$$

PCA - podiel vysvetlenej variancie

- $s_1 = 1$. hlavný komponent (PC)
 - vysvetľuje najväčší podiel variancie v dátach spomedzi všetkých normovaných lineárnych kombinácií pôvodných premenných
 - $s_k = k$ -ty hlavný komponent (PC), $k = 2, \dots, p$
 - vysvetľuje najväčšiu časť zostávajúcej variability v dátach, ktorá nebola vysvetlená pomocou s_1, \dots, s_{k-1} , spomedzi všetkých normovaných lineárnych kombinácií pôvodných premenných, ktoré sú nekorelované s prvými $(k - 1)$ hl. komponentami s_1, \dots, s_{k-1}
- počet hlavných komponentov je p , ale posledné PC vysvetľujú len malý podiel variability v dátach
- ⇒ posledné PC môžeme zanedbať
 - ⇒ pracujeme len s prvými $K \ll p$ PC

PCA - odhad hlavných komponentov I

1 odhad vektora \mathbf{u}

- $\hat{u}_j = \frac{1}{n} \sum_{l=1}^n x_{lj}$
- $\hat{\mathbf{u}} = \text{mean}(X)$

2 odhad kovariančnej matice

- $\hat{\Sigma}_{ij} = \frac{1}{n} \sum_{l=1}^n (x_{li} - \hat{u}_i)(x_{lj} - \hat{u}_j)$
- $\hat{\Sigma} = \text{cov}(X) \in \mathbb{R}^{p \times p}$

PCA - odhad hlavných komponentov II

- ③ stĺpce matice C = vlastné vektory matice $\hat{\Sigma}$

$$\hat{\Sigma}\hat{C} = \hat{C} \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_p \end{bmatrix}; \quad \lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \dots \lambda_{p-1} \geq \lambda_p$$

→ vlastné hodnoty súvisia s podielom vysvetlenej variancie v dátach

$$\alpha_i = \frac{\lambda_i}{\sum_{k=1}^p \lambda_k} \text{ je podiel variance vysvetlenej } i\text{-tym PC}$$

→ $[C, D] = \text{eig}(\hat{\Sigma}) \quad \lambda = \text{diag}(D)$

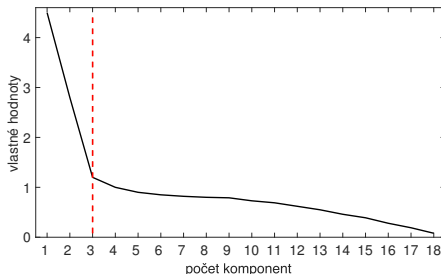
④ $\hat{S} = (X - \mathbf{1}_n \hat{\mathbf{u}}^T) \hat{C}$

PCA - určenie počtu hlavných komponentov K

- podiel vysvetlenej variancie \geq zvolený prah ρ , napr. $\rho = 0.95$

$$K \in \operatorname{argmin}_{k \in \{1, \dots, p\}} \frac{\sum_{i=1}^k \lambda_i}{\sum_{j=1}^p \lambda_j} \geq \rho$$

- lakťový graf (elbow diagram)



- Kaiserovo pravidlo

$$K \in \operatorname{argmax}_{k \in \{1, \dots, p\}} \lambda_k > \frac{1}{p} \sum_{i=1}^p \lambda_i$$

$$[C, S, \lambda, \sim, \text{expl}, u] = \text{pca}(X)$$

X = matica pozorovaní $n \times p$, MATLAB ju automaticky centruje

C = matica $p \times p$

- i -ty stĺpec reprezentuje i -ty hlavný komponent
- j -ty riadok reprezentuje príspevok j -tej pôvodnej premennej k jednotlivým hlavným komponentom

S = $n \times p$ matica hodnôt hl. komponentov pre n pozorovaní

λ = vlastné hodnoty patriace k jednotlivým hl. komponentom

$\text{expl} = (\alpha_1, \dots, \alpha_p)^T \times 100$, podiel vysvetlenej variancie jednotlivými hl. komponentami (v %)

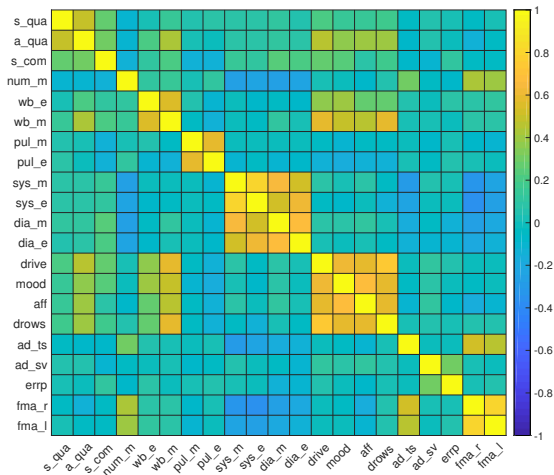
$u \in \mathbb{R}^p$ - priemer po stĺpcoch matice X

- [Rosipal et al., 2013]
- 148 ľudí, 2 noci v spánkovom laboratóriu → 296 pozorovaní
- dotazníky a testy ohľadom ich subjektívneho a objektívneho stavu (ráno, večer)

→ 21 premenných

- Self-rating questionnaire for sleep quality, awakening quality and somatic complaints
- Numerical memory test
- Well-being self-assessment scale evening/morning
- Pulse rate evening/morning
- Systolic and diastolic blood pressure evening/morning
- Visual analog scale test for drive, mood, affectivity and drowsiness
- Alphabetical cross-out test - total score, variability, % or errors
- Fine-motor activity test for right/left hand

PCA - Príklad: Denné miery



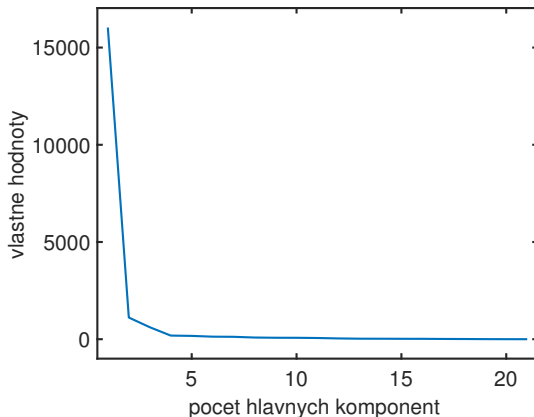
→ viaceré premenné sú korelované → nositelia podobnej informácie

→ $R = \text{corr}(\text{data})$

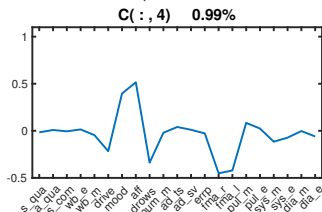
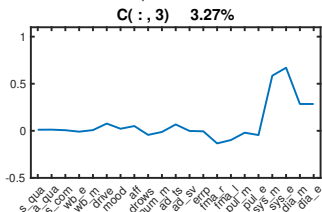
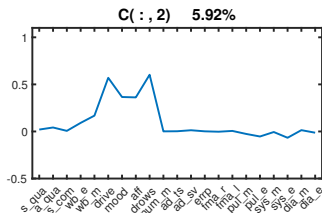
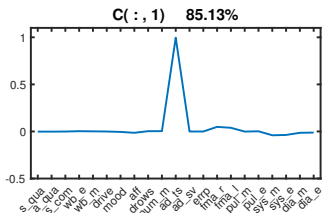
`heatmap(nazvy_premennych, nazvy_premennych, R)`

PCA - Príklad: Denné miery

- pôvodné dáta
- prvé 4 PC vysvetľujú viac ako 95% variability v dátach
- prvé 2 PC majú príslušné vlastné hodnoty vyššie ako priemer (Kaiserovo pravidlo)

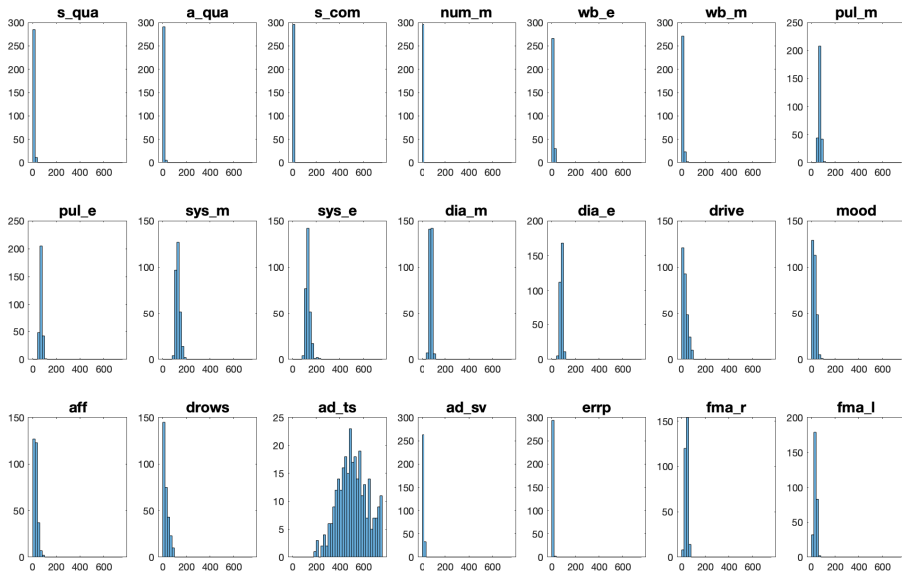


PCA - Príklad: Denné miery



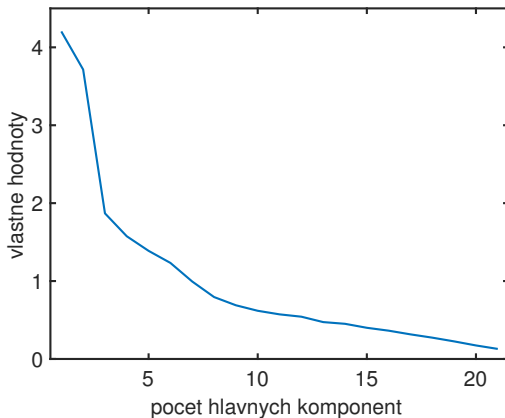
- C(:,1) - Abecedný pamäťový test → **výrazne vyššie hodnoty**
- C(:,2) - Visual analog scale for drive, mood, affectivity and drowsiness
- C(:,3) - krvný tlak

PCA - Príklad: Denné miery

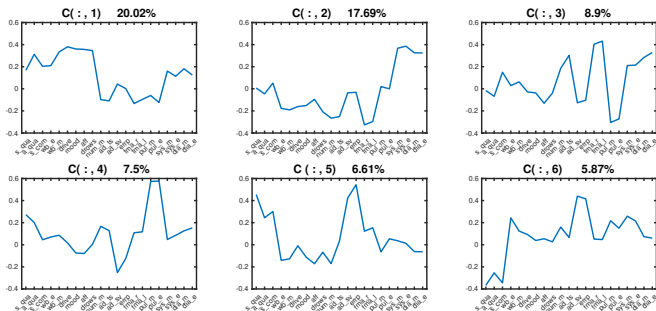


PCA - Príklad: Denné miery

- štandardizované dáta → $X_{new} = zscore(X)$
- prvých 17 PC vysvetľuje viac ako 95% variability v dátach
- prvých 6 PC majú príslušné vlastné hodnoty vyššie ako priemer (Kaiserovo pravidlo)



PCA - Príklad: Denné miery



- $C(:, 1)$ - subjektívny stav organizmu po prebudení
- $C(:, 2)$ - fyziologický stav organizmu (krvný tlak)
- $C(:, 3)$ - kognitívno-fyziologický stav organizmu (num. test, test jemnej motoriky, krvný tlak)
- $C(:, 4)$ - pulz ráno a večer

- **výhody**

- dochádza len k transformácií údajov, nie sú potrebné žiadne špeciálne predpoklady
- vhodné na vizualizáciu dát, resp. lepšie pochopenie charakteru dát
- odhadnuté PC možno použiť ako vstup v ďalších analýzach

- **nevýhody**

- nie je škálovo invariantná → **preškálovanie vedie k inému riešeniu**
- ak niektorá premenná vykazuje výrazne vyššie/nížšie hodnoty ako ostatné premenné
 - hl. komponenty budú v prvom rade zahŕňať informáciu z tejto premennej
- ⇒ niekedy je lepšie pracovať so **standardizovanými** premennými, t.j. zero-mean, unit-variance

Metóda hl. komponentov pre časové rady a funkcionálne dáta

Functional principal component analysis

PCA pre funkcionálne dáta a časové rady (FPCA)

- literatúra:
 - Ramsay, J. O. and Silverman, B. W. (2005). *Functional data analysis*. Springer Series in Statistics, New York, second edition.
 - Yao, F. et al. (2003). *Shrinkage estimation for functional principal component scores with application to the population kinetics of plasma folate*. Biometrics, 59(3):676-685.
- n pozorovaní (časových radov) x_1, \dots, x_n na intervale T

$$x_i(t) = \mu(t) + \sum_{k=1}^K \alpha_{ik} \varphi_k(t) + \varepsilon_i(t), \quad t \in T,$$

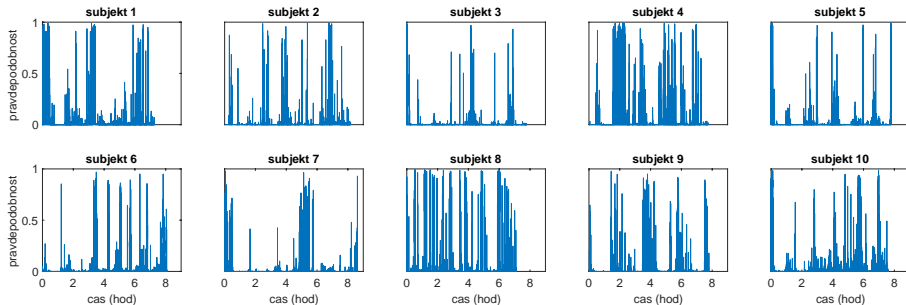
$$\int_T \varphi_k^2(t) = 1, \quad k = 1, \dots, K$$

$$\int_T \varphi_k(t) \varphi_l(t) = 0, \quad k \neq l$$

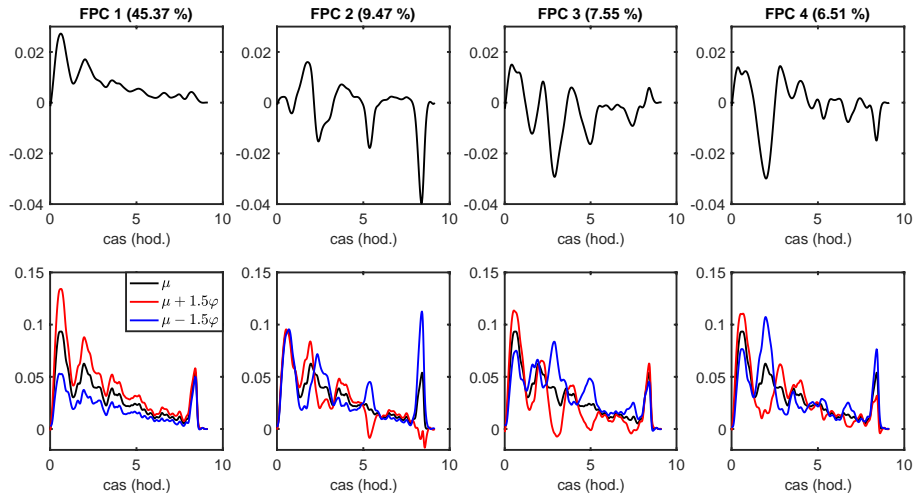
- nie je priamo implementovaná v MATLAB-e
- balík PACE - Principal Analysis by Conditional Estimation
 - <https://anson.ucdavis.edu/~mueller/data/pace.html>
 - Yao, F., Müller, H.G., Clifford, A.J., Dueker, S.R., Follett, J., Lin, Y., Buchholz, B., Vogel, J.S. (2003). Shrinkage estimation for functional principal component scores, with application to the population kinetics of plasma folate. *Biometrics* 59, 676-685.
 - funkcia *FPCA*

FPCA - Príklad: Spánkové krivky

- 146 ľudí
- 2 noci v spánkovom laboratóriu
- 292 kriviek reprezentujúcich stav bdelosti počas noci



FPCA - Příklad: Spánkové krivky



Faktorová analýza

Exploratory factor analysis (EFA)

- Spearman, C. (1904). *"General intelligence," objectively determined and measured*. The American Journal of Psychology, 15(2):201-292.
- **cieľ**: nájsť K nových priamo nepozorovateľných (tzv. skrytých, latentných) premenných, ktoré
 - sú navzájom nekorelované
 - popisujú dáta rovnako "kvalitne" ako pôvodné premenné
 - vieme interpretovať
- počet nových premenných K je vstupný parameter
- skryté (latentné) premenné nazývame aj faktory

$$X = \mathbf{1}_n \mathbf{u}^T + SC^T + E$$

$$x_{ij} = u_j + \sum_{k=1}^K s_{ik} c_{kj} + e_{ij}$$

→ $S \in \mathbb{R}^{n \times K}$

→ matica skrytých faktorov, ktoré sú navzájom nekorelované

→ $C \in \mathbb{R}^{p \times K}$

→ matica faktorových nákladov (váh)

→ prepojenie medzi pôvodnými premennými a skrytými faktormi

→ C, S sa odhadujú inak ako v PCA

- C - metóda hlavných faktorov; metóda maximálnej vierohodnosti; ...
- S - metóda regresnej analýzy, ...

EFA - odhad počtu faktorov

- apriórna vedomosť
- pokus-omyl s dôrazom na interpretáciu
- formálny postup
 - 1 korelačná matica dát $R = \text{corr}(X)$
 - 2 vypočítame tzv. redukovanú korelačnú maticu R^*
 - r je diagonála inverznej korelačnej matice R^{-1}
 - hlavnú diagonálu v R nahradíme $1 - \frac{1}{r}$
 - v MATLAB-e:

$$R^* = R - \text{eye}(p) + \text{diag}(1 - 1/\text{diag}(\text{inv}(R)))$$

- 3 vypočítame vlastné hodnoty matice R^*
 - v MATLAB-e: $e = \text{eig}(R^*)$
 - elbow diagram - hľadáme výrazný zlom v grafe
 - Kaiserovo pravidlo - $K = \#\{e_i : e_i > 1\}$

→ len ako odporúčané K , dôraz sa kladie na interpretáciu

EFA - rotácia faktorov

- **použitie:** pôvodné faktory sú ťažko interpretovateľné
- nech \mathbf{A} je rotačná matica a \mathbf{A}^{-1} reprezentuje inverznú rotáciu

$$X - \mathbf{1}_n \mathbf{u}^T = S C^T + E$$

$$X - \mathbf{1}_n \mathbf{u}^T = S A^{-1} A C^T + E$$

$$X - \mathbf{1}_n \mathbf{u}^T = (S A^{-1})(A C^T) + E$$

$$X - \mathbf{1}_n \mathbf{u}^T = S^* C^{*T} + E$$

- model s rovnakou chybou, ale **INO** interpretáciou faktorov
- môže zlepšiť interpretáciu faktorov

- najčastejšie používané rotácie:
 - *varimax* - maximalizuje varianciu štvorcov faktorových nákladov
 - *orthomax*, *quartimax*, *equamax*, *parsimax*
 - "vlastná"
 - ...

$[C, \sim, T, stats, S] = \text{factoran}(X, K, Name, Value)$

$X = n \times p$ matica pozorovaní alebo $p \times p$ kovariančná matica pozorovaní

K = počet skrytých (latentných) faktorov

• voliteľné parametre

- $Xtype = 'data', 'covariance'$
- $Rotate = \text{bez rotácie ('none')};$ s rotáciou ($'varimax', 'orthomax', \dots$)
- $Scores = \text{metóda na odhad } S ('wls', 'Bartlett', 'regression', \dots)$

$C = p \times K$ matica faktorových nákladov (váh)

T = rotačná matica

$stats$ = test hypotézy, že v dátach máme K skrytých faktorov

- $stats.p = p\text{-hodnota}$ ($stats.p < 0.05 \rightarrow$ hypotézu zamietame)

$S = n \times K$ matica skrytých faktorov

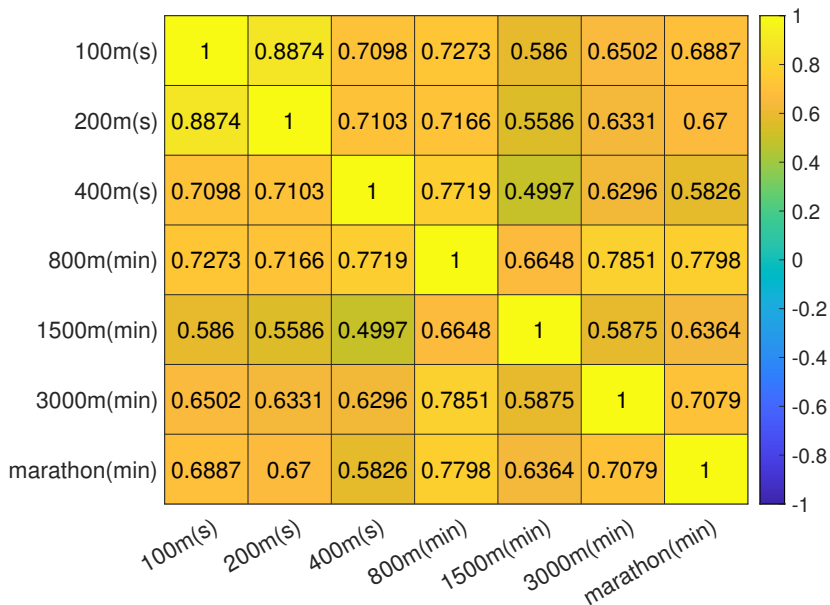
- podiel vysvetlenej variancie: $\alpha = \text{diag}(C^T C) / p$

- najlepšie výsledky reprezentantiek 55 krajín v 7 disciplínach

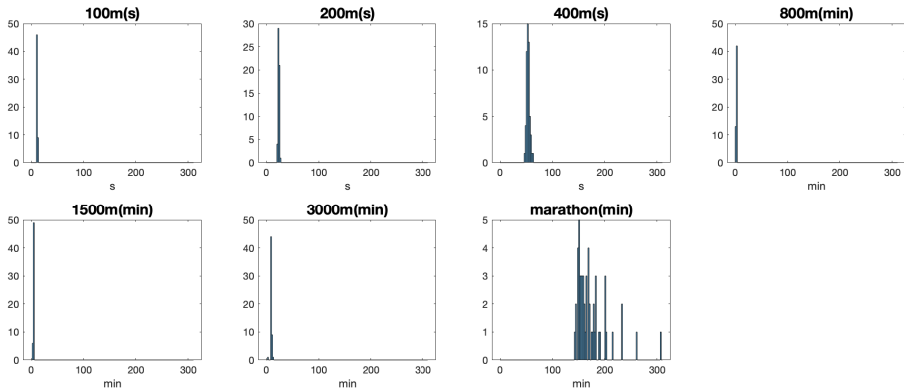
- 100 m (v s)
- 200 m (v s)
- 400 m (v s)
- 800 m (v min)
- 1500 m (v min)
- 3000 m (v min)
- maratón (v min)

→ chceme nájsť skryté faktory súhrnne charakterizujúce tieto disciplíny

EFA - Príklad: Výsledky bežeckých disciplín u žien



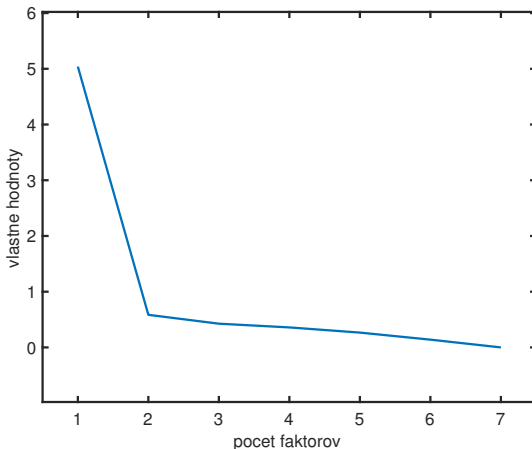
EFA - Príklad: Výsledky bežeckých disciplín u žien



→ potrebná štandardizácia dát $X^* = zscore(X)$

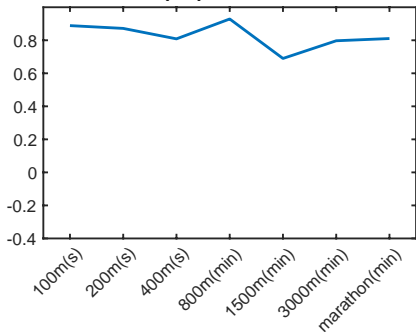
EFA - Príklad: Výsledky bežeckých disciplín u žien

- stanovenie počtu faktorov
 - Kaiserovo pravidlo: 1 faktor
 - elbow diagram: 1-2 faktory

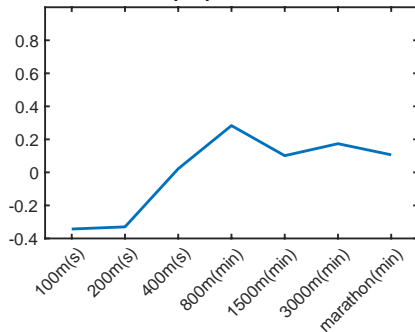


EFA - Príklad: Výsledky bežeckých disciplín u žien

C(:,1) - 69.07 %

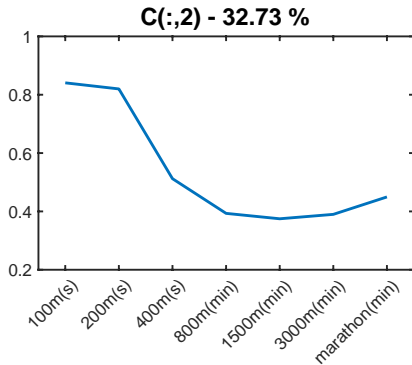
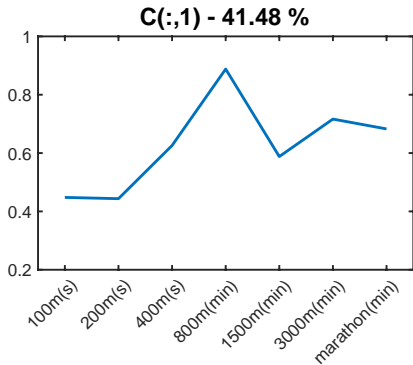


C(:,2) - 5.13 %



- model s 2 faktormi, **bez rotácie faktorov**
 - C(:,1) - (+) pre všetky disciplíny (rovnomé váhy)
→ faktor by mal reprezentovať len malú skupinu premenných
 - C(:,2) - (+) pre vytrvalostné disciplíny, (-) pre rýchlostné disciplíny

EFA - Príklad: Výsledky bežeckých disciplín u žien



- model s 2 faktormi, *varimax rotácia faktorov*
 - C(:,1) - vyššie (+) hodnoty pre vytrvalostné disciplíny
 - C(:,2) - vyššie (+) hodnoty pre rýchlostné disciplíny

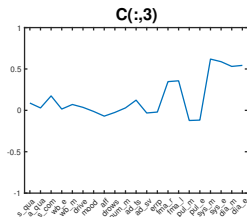
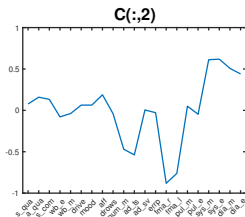
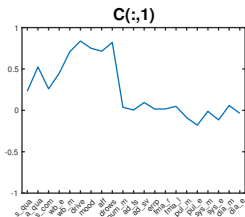
- [Rosipal et al., 2013]
- 148 ľudí, 2 noci v spánkovom laboratóriu → 296 pozorovaní
- dotazníky a testy ohľadom ich subjektívneho a objektívneho stavu (ráno, večer)

→ 21 premenných

- Self-rating questionnaire for sleep quality, awakening quality and somatic complaints
- Numerical memory test
- Well-being self-assessment scale evening/morning
- Pulse rate evening/morning
- Systolic and diastolic blood pressure evening/morning
- Visual analog scale test for drive, mood, affectivity and drowsiness
- Alphabetical cross-out test - total score, variability, % or errors
- Fine-motor activity test for right/left hand

EFA - Príklad: Denné miery

- štandardizácia dát $X_{new} = zscore(X)$
- EFA s 3 faktormi, **bez rotácie** faktorov

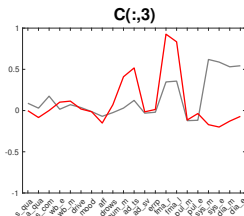
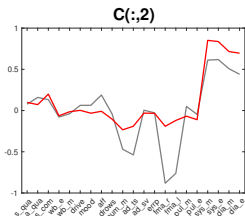
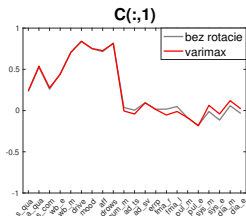


- $C(:, 1)$ (17,23%) - faktor subjektívneho hodnotenia kvality spánku
- $C(:, 2)$ (15,15%) - krvný tlak (+) + “niektoré” kognitívne testy (-)
- $C(:, 3)$ (7,87%) - krvný tlak (+) + “zvyšné” kognitívne testy (-)

→ **problém:** pôvodné premenné by mali mať “vysoké” (+/-) váhy len pre jeden z latentných faktorov

EFA - Príklad: Denné miery

- štandardizácia dát $X_{new} = zscore(X)$
- výsledky v [Rosipal et al., 2013]
- EFA s 3 faktormi, **rotácia** faktorov pomocou metódy *varimax*



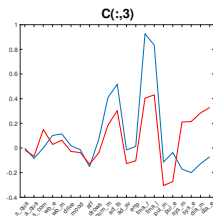
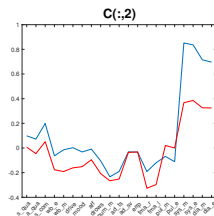
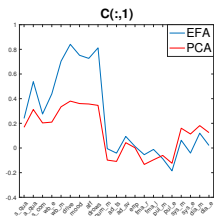
- $C(:, 1)$ (17,36%) - faktor subjektívneho hodnotenia kvality spánku
- $C(:, 2)$ (12,63%) - faktor fyziologického stavu organizmu
- $C(:, 3)$ (10,25%) - faktor kognitívneho stavu organizmu po prebudení

EFA

- predpoklad existencie latentných premenných
- K je vstupný parameter
- odhady - metóda hlavných faktorov, metóda maximálnej vierohodnosti

PCA

- nové premenné = lin. kombináciu pôvodných
 - K sa určí až po analýze
 - odhady - eigenanalysis, SVD, ...
- výstup PCA = inicializačné odhady faktorov v EFA





Rosipal, R., Lewandowski, A., and Dorffner, G. (2013).

In search of objective components for sleep quality indexing in normal sleep.
Biological Psychology, 94(1):210–220.