

CALCUL NUMERIC

Analiza Componentelor Principale

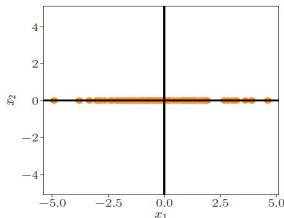
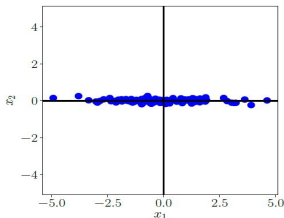
Paul Irofti
Cristian Rusu
Andrei Pătrașcu

Departament Informatică
Universitatea din București

- **Reducție Dimensională**
- Algoritm PCA
- Aplicație

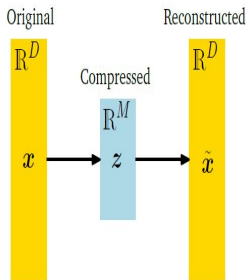


- Utilizarea și prelucrarea datelor de dimensiuni mari impune adesea mari dificultăți privind analiza, vizualizarea și stocarea acestora
- Cu toate acestea, dimensiunile mari atrag redundanțe ascunse în date, e.g. anumite componente pot fi exprimate ca și combinații liniare de alte componente.



Analiza componentelor principale (PCA) reprezintă o metodă simplă de reducere dimensională liniară, apărută pentru prima oară în lucrările lui Pearson (1901)

- Datele: $X = [x_1, x_2, \dots, x_N] \in \mathbb{R}^{d \times N}$ cu media 0
- PCA determină o bază subspațiului optim de dimensiune $m \ll d$, i.e. $B \in \mathbb{R}^{d \times m}$, care reține maximum de informație conținută în X
- Astfel, reduce dimensiunea de la d la m , $X \rightarrow Z$

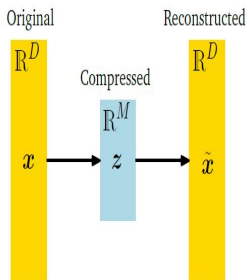


Matricea bazei $B = [b_1 \ b_2 \ \cdots \ b_m]$, unde $m < d$, este ortonormală:

- $b_i^T b_j = 0$ pentru orice $i \neq j$
- $b_i^T b_i = 1$ pentru orice $1 \leq i \leq N$

Coordonatele vectorilor în noua bază: $z_i = B^T x_i$

Proiecțiile vectorilor în noua bază: $\hat{x}_i = Bz_i$



Subspațiul optim corespunde bazei B care minimizează distanța $\|\hat{x}_i - x_i\|$.
Echivalent, baza care maximizează varianța noilor coduri $\{z_i\}_{i=1}^N$ (Hotelling)

$$V_z[z] = V_x[B^T(x - \bar{x})] = V_x[B^T x - B^T \bar{x}] = V_x[B^T x],$$

unde am folosit media $\bar{x} = 0$.

Reamintim: varianța unei variabile aleatoare discrete y (cu N valori) este

$$V_y[y] = \frac{1}{N} \sum_i (y_i - \bar{y})^2$$

unde \bar{y} este media lui y .



- Reducție Dimensională
- **Algoritm PCA**
- Aplicație



Începem prin maximizarea varianței pentru prima componentă din z :

$$\begin{aligned} V[z_i] &= \frac{1}{N} \sum_{i=1}^N z_{i1}^2 \\ &= \frac{1}{N} \sum_{i=1}^N (b_1^T x_i)^2 = \frac{1}{N} \sum_{i=1}^N b_1^T (x_i x_i^T) b_1 \\ &= b_1^T \underbrace{\left(\frac{1}{N} \sum_{i=1}^N x_i x_i^T \right)}_S b_1 \end{aligned}$$

Notăm: $S := \frac{1}{N} \sum_{i=1}^N x_i x_i^T$ matricea de covarianță a datelor X



Prima iterație necesită rezolvarea:

$$\begin{aligned} \max \quad & b_1^T S b_1 \\ \text{s.l.} \quad & \|b_1\| = 1 \end{aligned}$$

Condiții de optimalitate: punctele de maxim satisfac setul de egalități

$$\begin{aligned} S b_1 &= \lambda_1 b_1 \\ \|b_1\| &= 1 \end{aligned}$$

- Recunoaștem aceste relații din cursul de valori/vectori proprii
- Soluția: (b_1, λ_1) vectorul propriu maximal al lui S
- b_1 reprezintă prima componentă principală
- λ_1 varianța datelor proiectate pe subspațiul definit de prima componentă principală



Considerăm că au fost executate $k - 1$ iterații, care au produs $\{b_1, \dots, b_{k-1}\}$
 Mai departe, la iterația k , eliminăm "efectul" primelor $k - 1$ componente

$$\hat{X}^{(k)} := X - \sum_{i=1}^{k-1} b_i b_i^T X \quad \hat{S}^{(k)} := \frac{1}{N} \hat{X}^{(k)} (\hat{X}^{(k)})^T$$

Aplicăm iterația k pe noile date proiectate $\hat{X}^{(k)}$:

$$\begin{aligned} \max \quad & b_k^T \hat{S}^{(k)} b_k \\ \text{s.l.} \quad & \|b_k\| = 1 \end{aligned}$$

Vectorii proprii ai lui $\hat{S}^{(k)}$ sunt identici cu cei ai lui S



Vectorii proprii ai lui $\hat{S}^{(k)}$ sunt identici cu cei ai lui S

$$\begin{aligned}\hat{S}^{(k)} b_i &= \frac{1}{N} \hat{X}^{(k)} (\hat{X}^{(k)})^T b_i = \frac{1}{N} \left(X - \sum_{i=1}^{k-1} b_i b_i^T X \right) \left(X - \sum_{i=1}^{k-1} b_i b_i^T X \right)^T b_i \\ &= \begin{cases} 0, & \text{dacă } i < k \\ \lambda_i b_i, & \text{dacă } i \geq k \end{cases}\end{aligned}$$

Verificați!



Din relațiile precedente, avem $\lambda_i(\hat{S}^{(k)}) = \begin{cases} 0, & \text{dacă } i < k \\ \lambda_i(S), & \text{dacă } i \geq k \end{cases}$

De aceea, subproblema asociată iterației k :

$$\begin{aligned} \max \quad & b_k^T \hat{S}^{(k)} b_k \\ \text{s.l.} \quad & \|b_k\| = 1 \end{aligned}$$

are soluția: *vectorul propriu asociat valorii proprii maxime a matricii $\hat{S}^{(k)}$*

Concluzie: Baza optimă PCA, i.e. $B = [b_1 \ b_2 \ \cdots \ b_m]$, este formată din vectorii proprii ai matricii S asociați celor mai mari m valori proprii.



Informația conservată prin reconstrucție pe baza lui B este data de:

$$V_m := \sum_{i=1}^m \lambda_i(S)$$

Informația pierdută (eroarea) la reconstrucție:

$$J_m := \sum_{i=m+1}^d \lambda_i(S)$$



Pași PCA:

- 1 Se compune $S = \frac{1}{N}XX^T$, unde $X = [x_1 \cdots x_N]$
- 2 Calculează DVP (Descompunerea Valorilor Proprii) lui $S = U\Lambda U^T$
- 2' Calculează DVS (Descompunerea Valorilor Singulare) lui $X = U\Sigma V^T$
- 3 Calcul subspațiu: $B = [u_1 \cdots u_m]$
- 4 Reconstrucție: $Z = B^T X$; $\tilde{X} = BZ$

Pentru pasul al doilea se execută: 2 SAU 2'

- Coloanele lui U reprezintă vectorii proprii ai matricii XX^T
- Între valorile proprii λ_i (pas 2) și valorile singulare (pas 2') are loc relația:

$$\lambda_i = \frac{\sigma_i^2}{N}$$



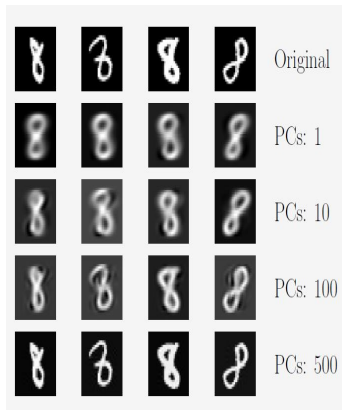
- Reducție Dimensională
- Algoritm PCA
- **Aplicație**



Aplicație MNIST:

- Considerăm setul de date MNIST format din imagini ale cifrelor 0-9 scrise de mână
- Se echivalează fiecare pixel cu un număr egal cu nivelul de gri și se vectorizează
- După transformarea imaginilor în vectori numerici avem:
 $N = 5389, d = 784, m \in \{1, 10, 100, 500\}$
- Numărul maxim de componente principale: 784





Concluzii aplicație MNIST:

- Pe prima linie apar 4 imagini originale ale cifrei 8
- Pe fiecare dintre liniile următoare transformările în imagini ai vectorilor reconstruiți folosind un număr fixat de componente principale
- La final de linie este indicat numărul de componente principale folosit
- Observăm că reconstrucția pe baza primelor 500 de componente este aproape perfectă
- **Compresie:** stocăm doar primele 500 CP (< 784) + coordonatele asociate imaginilor $z_i = B^T x_i$
- Necesari memorie:
$$\underbrace{500 \cdot 784}_{\text{Mem. CP}} + \underbrace{500 \cdot N}_{\text{Mem. coduri } z} < \underbrace{784 \cdot N}_{\text{Mem. imagini originale}}$$



- Deisenroth, M. P., Faisal, A. A., & Ong, C. S. (2020). Mathematics for machine learning. Cambridge University Press.

