

11. Analiza componentelor principale

Normalizarea datelor
Baze ortogonale de reprezentare a datelor
Maximizarea varianței datelor

11.1 Motivație

Obiectivele algoritmului de *Analiză a Componentelor Principale* (eng. Principal Components Analysis - PCA) sunt analiza datelor în scopul identificării de structuri în date și de reducere a numărului de dimensiuni prin care sunt reprezentate datele. Algoritmul schimbă baza de reprezentare a datelor, utilizând primele componente principale ce descriu datele și ignorând restul. Primul component principal este direcția care maximizează varianța datelor proiectate în noua bază de reprezentare.

Formal, datele sunt reprezentate într-un spațiu n -dimensional: $x \in \mathbb{R}^n$. PCA identifică un subspațiu k -dimensional al datelor (unde $k < n$) prin calculul vectorilor proprii ai lui x .

Considerăm un set de date de antrenare:

$$\{x^{(i)}; i = 1, \dots, m\} \quad (11.1)$$

ce conține atribute (precum viteza maximă, comportamentul în viraje, etc) pentru m tipuri de automobile.

Fie $x \in \mathbb{R}^n$ pentru fiecare exemplu de antrenare i ($n \ll m$). Fără a ne fi cunoscut, două tipuri de atribute (x_i , respectiv x_j) reprezintă:

1. x_i : viteza maximă a automobilului în kilometri pe oră
2. x_j : viteza maximă a automobilului în mile pe oră.

Aceste două atribute sunt liniar dependente. Astfel, datele sunt, de fapt, reprezentate într-un spațiu $n - 1$ dimensional. Metoda PCA este utilizată în identificarea și înlăturarea acestor redundanțe.

Un alt exemplu este acela al unei baze de date compusă din părerile unor piloți de drone controlate prin telecomandă, unde:

1. $x_1^{(i)}$ este o mărime a gradului de îndemânare al pilotului i ,

2. $x_2^{(i)}$ reprezintă gradul de bucurie pe care îl are în pilotare.

Deoarece dronele controlate prin telecomandă sunt dificil de manevrat, doar piloții dedicați, care au o pasiune pentru această activitate, ajung să fie buni piloți. Astfel, atributele x_1 și x_2 sunt puternic corelate. După cum este redat în Figura 11.1, se poate spune că datele sunt reprezentate de-a lungul unei axe diagonale (direcția u_1) ce descrie modul de comportare al unui pilot. În cele ce urmează vom calcula direcția u_1 prin metoda PCA.

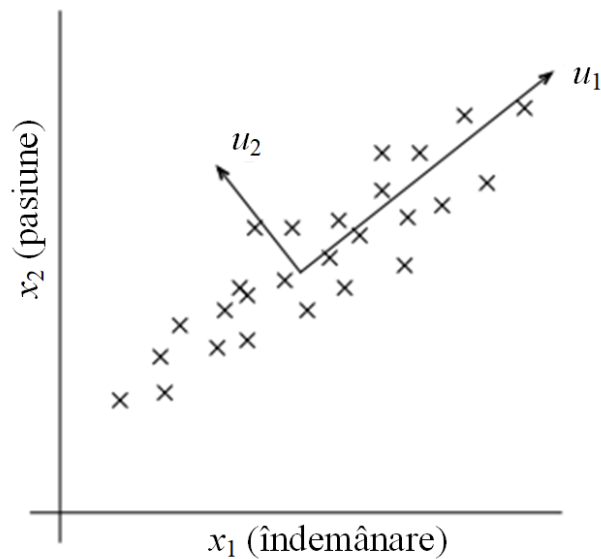


Fig. 11.1 Distribuția caracteristicilor ce descriu comportamentul unui pilot de dronă.

11.2 Vectorii Proprii ai unei Matrici

Valorile proprii reprezintă un set de variabile scalare asociat unui sistem liniar de ecuații exprimat sub formă matriceală. Acestea sunt întâlnite și sub denumirea de rădăcini caracteristice, valori caracteristice, valori proprii, sau rădăcini latente.

Fiecărei valori proprii îi corespunde un vector propriu. Vectorul propriu al unei transformări liniare pe un spațiu vectorial este un vector nenul a cărui direcție rămâne neschimbată de către acea transformare.

Fie Σ o matrice pătratică $\mathbb{R}^{m \times m}$ ce reprezintă o transformare liniară. Un vector u nenul de dimensiune $m \times 1$ (vector coloană) se numește vector propriu pentru transformarea liniară Σ , dacă există un scalar λ , astfel încât:

$$\Sigma \mathbf{u} = \lambda \mathbf{u} \quad (11.2)$$

unde u este vectorul propriu corespunzător valorii proprii λ . Vectorii și valorile proprii se pot obține aplicând metoda descompunerii în valori singulare (Singular Value Decomposition - SVD) a unei matrici.

11.3 Normalizarea Datelor

Anterior dezvoltării algoritmului PCA, vom normaliza datele utilizând media și varianța lor, astfel:

1. Fie $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$;
2. Fiecare $x^{(i)}$ va fi înlocuit cu $x^{(i)} - \mu$;
3. Fie $\sigma_j^2 = \frac{1}{m} \sum_i \left(x_j^{(i)}\right)^2$;
4. Fiecare $x_j^{(i)}$ va fi înlocuit cu $x_j^{(i)} / \sigma_j$.

Pașii (1) și (2) centrează datele pe media zero. Acești pași pot fi omiși atunci când se cunoaște faptul că datele au media zero (spre exemplu, pentru seriile de timp ce descriu semnale acustice).

Pașii (3) și (4) scalează fiecare coordonată a datelor la varianța unitate. Această scalare asigură că diferitele atribute au aceeași scală și sunt astfel tratate la fel. Spre exemplu, dacă x_1 este viteza maximă a unui autoturism, măsurată în km/h (ce are valori de ordinul a sute de km/h), iar x_2 reprezintă numărul de scaune (ce în mod normal sunt în număr de 2 sau 4), atunci scalarea prin varianță face atributele mai compatibile între ele.

Pașii (3) și (4) pot fi omiși atunci când știm că atributele sunt reprezentate la o scală compatibilă. Un exemplu în acest sens este atunci când considerăm fiecare pixel dintr-o imagine gri ca fiind o caracteristică $x_j^{(i)}$, fiecare valoare variind în intervalul $[0, 255]$, ce corespunde intensității pixelului j din imaginea i .

11.4 Algoritmul PCA

Componenta principală i este direcția ortogonală către primele $i - 1$ componente principale. Componentele principale sunt vectorii proprii ai matricei de covarianță a datelor.

Folosind datele normalizate, se calculează axa principală de variație u pe care sunt distribuite datele. O modalitate de descriere a acestei probleme este de a găsi vectorul unitate u , în așa fel încât varianța datelor proiectate să fie maximă atunci când sunt proiectate de-a lungul direcției lui u . Intuitiv, datele vor avea pentru început un anumit grad de varianță. Dorim să găsim o direcție a lui u în așa fel încât să se păstreze cât mai multă varianță atunci când datele sunt proiectate pe direcția/sub-spațiul lui u .

Considerați setul de date normalizat reprezentat în Figura 11.2.

O posibilă direcție pentru u este ilustrată în Figura 11.3. Cercurile reprezintă proiecția datelor originale pe direcția lui u .

Se poate observa din Figura 11.3 că datele proiectate păstrează încă un grad ridicat de varianță, cu puncte îndepărtate de valoarea zero. În contrast cu acest caz, varianța datelor

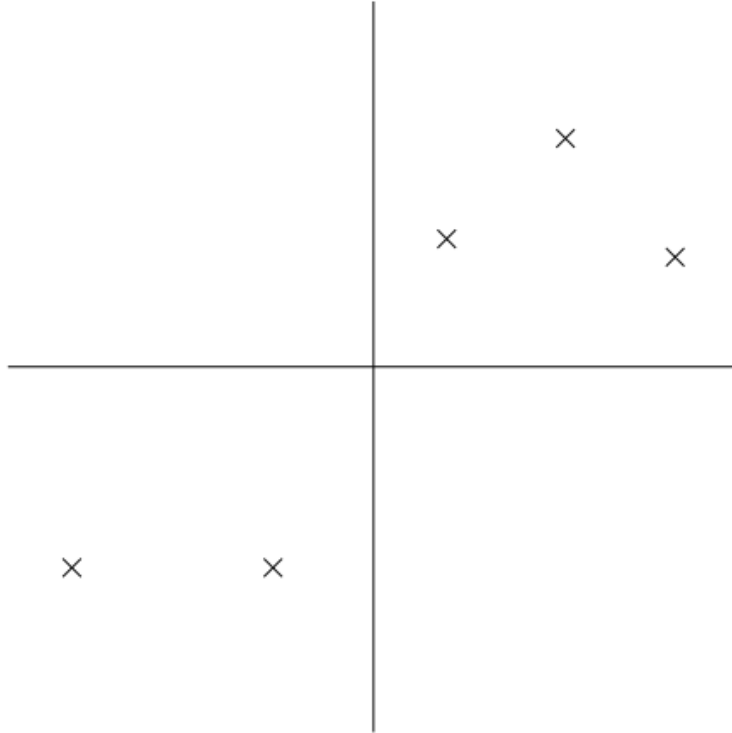


Fig. 11.2 Set de date normalizat.

proiectate de-a lungul direcției ilustrate în Figura 11.4 este mult mai mică, cu puncte mult mai apropiate de origine.

În analiza componentelor principale, dorim să selectăm automat direcția lui u ce corespunde Figurii 11.3. Formal, **pentru un vector unitate u și un punct x , lungimea proiecției lui x pe u este dată de $x^T u$** . Spre exemplu, dacă $x^{(i)}$ este un punct din baza de date originală, atunci proiecția sa pe u este distanța $x^T u$ de la origine. Astfel, pentru a maximiza varianța proiecțiilor, dorim să alegem vectorul unitate u în așa fel încât să maximizăm:

$$\begin{aligned} \frac{1}{m} \sum_{i=1}^m \left(x^{(i)T} u \right)^2 &= \frac{1}{m} \sum_{i=1}^m u^T x^{(i)} x^{(i)T} u \\ &= u^T \left(\frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T} \right) u. \end{aligned}$$

Se poate observa că maximizarea expresiei de mai sus utilizând constrângerea $\|u\|_2 = 1$ rezultă în vectorii proprii principali ai expresiei:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)T}. \quad (11.3)$$

Expresia 11.3 reprezintă matricea de covarianță a datelor, luând în considerare că datele au media zero (sunt centrate pe valoarea zero).

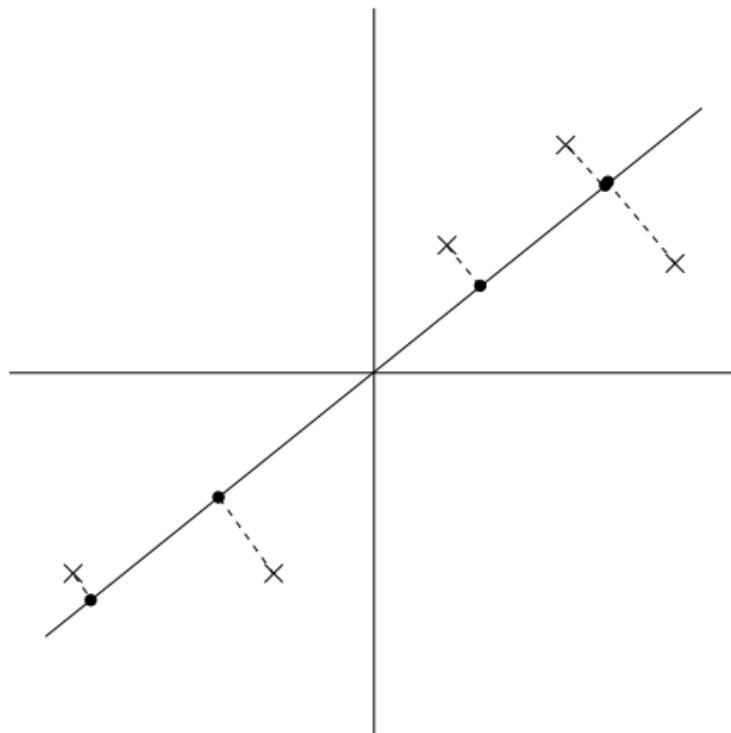


Fig. 11.3 Posibila proiecție a setului de date din Figura 11.2.

Astfel, pentru a găsi un subspațiu 1-dimensional ce poate aproxima datele 2-dimensionale, ar trebui ales u ca fiind vectorii proprii ai matricei Σ . Generalizat, dacă dorim să proiectăm datele într-un subspațiu k -dimensional ($k < n$), ar trebui să alegem u_1, u_2, \dots, u_k ca și primii k vectori proprii ai matricei Σ . Noile valori u_i reprezintă noua bază ortogonală a datelor. Deoarece Σ este simetrică, vectorii u_i pot fi aleși ortogonali unul față de celălalt.

Pentru a reprezenta $x^{(i)}$ în această nouă bază, trebuie să calculăm vectorii corespunzători:

$$y^{(i)} = \begin{bmatrix} u_1^T x^{(i)} \\ u_2^T x^{(i)} \\ \vdots \\ u_k^T x^{(i)} \end{bmatrix}, \quad y^{(i)} \in \mathbb{R}^k. \quad (11.4)$$

Luând în considerare că $x^{(i)} \in \mathbb{R}^n$, vectorul $y^{(i)}$ oferă o aproximare/reprezentare a lui $x^{(i)}$ într-un spațiu k -dimensional. Din această cauză metoda PCA este adesea întâlnită sub denumirea de algoritm de *reducere a dimensiunii*. Vectorii u_1, u_2, \dots, u_k sunt denumiți primele k componente principale ale datelor.

Cu toate că în acest curs vom reduce datele la subspațiul 1-dimensional ($k = 1$), se poate demonstra că dintre toate bazele ortogonale u_1, u_2, \dots, u_k , cea aleasă maximizează expresia $\sum_i \|y^{(i)}\|_2^2$. Astfel, alegerea unei baze ortogonale păstrează pe cât de mult posibil variabilitatea din datele originale.

Algoritmul PCA se utilizează într-o gamă largă de probleme, precum compresia datelor,

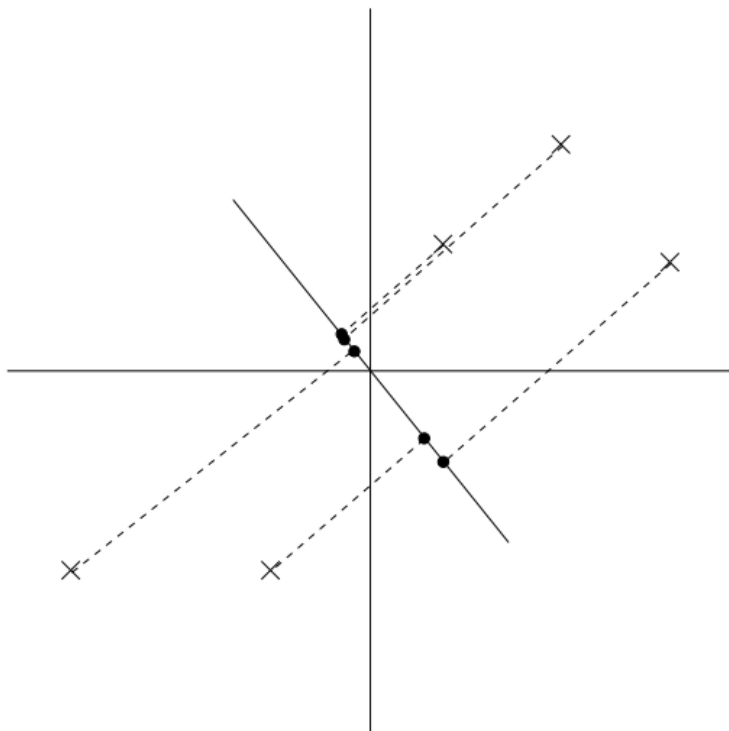


Fig. 11.4 Posibila proiecție a setului de date din Figura 11.2.

vizualizarea distribuției caracteristicilor (eng. features) într-un spațiu 1-, 2-, sau 3-dimensional, sau ca și o modalitate de reducere a zgomotului.

Bibliografie

- [1] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006.
- [2] I. Goodfellow, Y. Bengio and A. Courville, *Deep Learning*. MIT Press, 2016, <http://www.deeplearningbook.org>.
- [3] M. Lutz, *Learning Python*, 2nd Ed. Sebastopol, CA, USA: O'Reilly & Associates, Inc., 2003.
- [4] A. Ng, "Stanford cs229 - machine learning," 2008.
- [5] S. J. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, 2nd Ed. Pearson Education, 2003.
- [6] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser and I. Polosukhin, "Attention is all you need," *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc., 2017.