

Principal Component Analysis

Machine Learning 2022-2023 - UMONS

Souhaib Ben Taieb

Exercise 1

Consider the following design matrix X :

$$X = \begin{array}{cc} & \begin{array}{c} X_1 \quad X_2 \end{array} \\ \begin{array}{c} 4 \\ 2 \\ 5 \\ 1 \end{array} & \begin{array}{c} 1 \\ 3 \\ 4 \\ 0 \end{array} \end{array}$$

We want to represent the data in only one dimension using principal components analysis (PCA). To this end :

- Center the data.
- Compute the sample covariance matrix C .
- Compute the eigenvalues and eigenvectors of the covariance matrix C .
- Plot the dataset, and draw the first principal component direction (as a line) and the projections of all four sample points onto the principal direction.
- Label each data point with its principal component score.
- Compute the proportion of variance explained by the first principal component.
- Add the projections of the data points onto the second principal component, compute the second principal component scores, and show that the sum of the variance explained by each component is equal to the total variance of the data.

Recall that the eigenvalues of a square matrix A are obtained by solving the characteristic equation $\det(A - \lambda I) = 0$. If $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$, then $\det(A) = ad - bc$. Finally, an eigenvector v of A must satisfy $Av = \lambda v$. If X is centered, then its covariance matrix is equal to $C = \frac{1}{n}X^T X$.

Exercise 2

Suppose that the columns of $X \in \mathbb{R}^{n \times p}$ have been centered (i.e. they have sample mean zero). The total sample variance of X is defined as $\text{TV} = \sum_{j=1}^p \text{Var}(X_j) = \frac{1}{n} \text{Trace}(X^\top X)$.

Let $X = UDV^\top$ be the singular value decomposition of X where the columns of $U \in \mathbb{R}^{n \times p}$ and $V \in \mathbb{R}^{p \times p}$ are orthonormal and the matrix $D \in \mathbb{R}^{p \times p}$ is diagonal with positive real entries, $D = \text{diag}(d_1, \dots, d_n)$, with $d_1 \geq d_2 \geq \dots \geq d_n \geq 0$.

Prove that the total variance of X is given by $\frac{1}{n} \sum_i^n d_i^2$. How do you link this result to the eigenvalues of covariance matrix C ?

Hints :

- $\text{Trace}(AB) = \text{Trace}(BA)$
- $(AB)^\top = B^\top A^\top$
- If V is orthonormal, then $V^\top V = I$
- If A is diagonal, then $A^\top A = A^2$