

# 2.1 Data analytics for dimensionality reduction: Principal Component Analysis (PCA)

Prof. Massimiliano Grosso

University of Cagliari, Italy

massimiliano.grosso@dimcm.unica.it

GRICU PhD School 2021

*Digitalization Tools for the Chemical and Process Industries*

March 12, 2021

1



## Outline

- Motivations
- Basic concepts
- Preprocessing
- Mathematical background
- Dimension reduction
- Geometrical interpretation



2

1

## Motivations

- Concerns when dealing with “**huge**” amount of data:
  - The **size** of the data:
    - The **useful** information is often «**hidden**» amongst **hundred/thousands of variables**
  - The measurements are often **highly correlated** with one another (multicollinearity)
    - The number of **independent variables (degrees of freedom)** is **much less** than the number of measurements on hand
- **Noise** in the measurements
  - Difficulties in distinguishing the noise from the deterministic variations induced by external sources



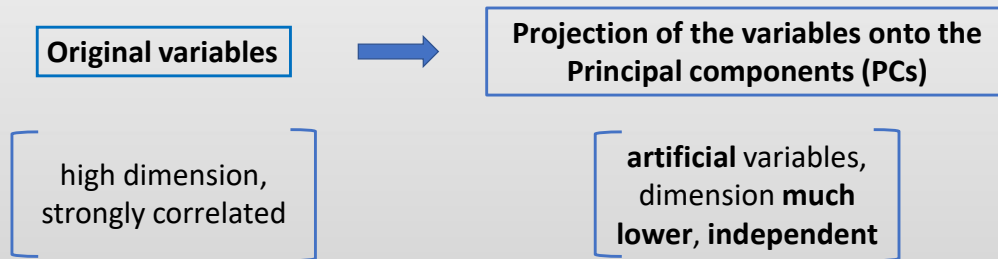
## Motivations

- Multivariate data analysis method for
  - Explorative data analysis
  - Outlier detection
  - Rank reduction
  - Graphical clustering
  - Classification
- PCA allows interpretation based on **all** variables simultaneously, leading to understanding deeper than what is possible looking at the individual variables alone
- It is the first **multivariate analysis** to be carried out



## PCA: Basic concepts

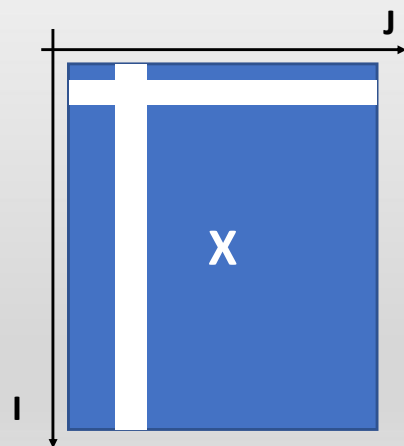
- Aim of the PCA:



5

## PCA: Basic concepts

- Data must be collected on matrix **X**
- Column vectors represent the **variables** ( $j=1, \dots, J$ )
  - attributes, wavelenghts, physical/chemical parameters etc.
- Row vectors represent the **samples** ( $i=1, \dots, I$ ) collected during the experiments



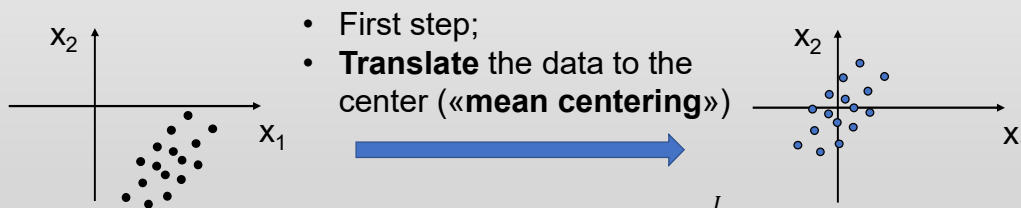
2.1 Principal Component Analysis (M. Grosso)

6

6

## Preprocessing of the data

- Matrix  $\mathbf{X}$  can be visualized in a coordinate system made up by
  - J orthogonal axes, each representing one of the original J variables
  - Each i-th sample is a J-dimensional row vector
- Two-dimensional example with two variables highly correlated



- First step;
- **Translate** the data to the center («mean centering»)

$$x_{ij} = x_{ij}^* - \bar{x}_{\bullet j}^* \quad \text{where} \quad \bar{x}_{\bullet j}^* = \frac{1}{I} \sum_{i=1}^I x_{ij}^*$$

## Preprocessing of the data

- Mean centering allows to consider the covariance matrix  $\mathbf{C} = \mathbf{X}^T \mathbf{X}$

$$\mathbf{X} = \begin{bmatrix} x_{11} & \cdots & x_{1j} \\ \vdots & \ddots & \vdots \\ x_{I1} & \cdots & x_{Ij} \end{bmatrix} = \begin{bmatrix} x_{11}^* - \bar{x}_{\bullet 1}^* & \cdots & x_{1j}^* - \bar{x}_{\bullet j}^* \\ \vdots & \ddots & \vdots \\ x_{I1}^* - \bar{x}_{\bullet 1}^* & \cdots & x_{Ij}^* - \bar{x}_{\bullet j}^* \end{bmatrix}$$

- Indeed, for the element  $kl$

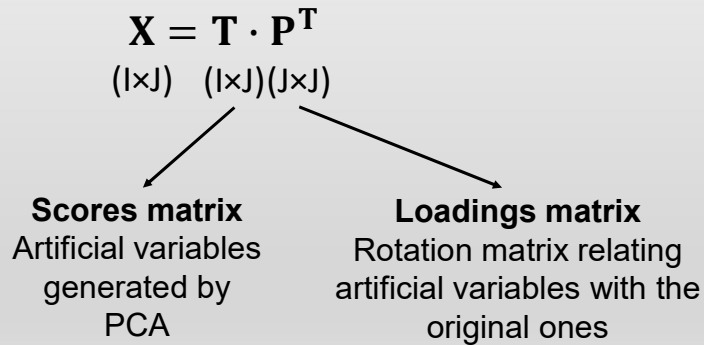
$$C_{kl} = (\mathbf{X}^T \mathbf{X})_{kl} = \sum_{i=1}^I (x_{ik}^* - \bar{x}_{\bullet k}^*)(x_{il}^* - \bar{x}_{\bullet l}^*) = \sum_{i=1}^I x_{ik} x_{il}$$

- The diagonal elements of  $\mathbf{C}$  are the *dispersion* related to the j-th variable

$$C_{jj} = (\mathbf{X}^T \mathbf{X})_{jj} = \sum_{i=1}^I (x_{ij}^* - \bar{x}_{\bullet j}^*)^2 = \sum_{i=1}^I x_{ij}^2$$

## PCA – Basic concepts

- Principal Component Analysis is based on the **decomposition** of the dataset matrix **X**



## PCA – Basic concepts

- Important properties:

1. Even the scores are mean centered

$$\bar{x}_{\cdot j} = 0 \quad \forall j = 1, \dots, J \quad \Rightarrow \quad \bar{t}_{\cdot j} = 0 \quad \forall j = 1, \dots, J$$

2. **Column vectors of the score matrix T are orthogonal:**

$$\mathbf{t}_m^T \mathbf{t}_n = 0 \quad \forall m \neq n$$

1. The square of the score matrix  $\mathbf{\Lambda} = \mathbf{T}^T \cdot \mathbf{T}$  is **diagonal**

3. **Loadings matrix P is orthogonal:**  $\mathbf{P}^T = \mathbf{P}^{-1} \Rightarrow \mathbf{P}^T \cdot \mathbf{P} = \mathbf{I}$

## Mathematical background

- PCA scores and loadings can be related to the computation of the **eigenvalues** and **eigenvectors** of the  $J \times J$  **covariance** matrix

$$\mathbf{C} = \mathbf{X}^T \mathbf{X}$$

- **Remark**
- **C** is a square, symmetric matrix, this leads to the following properties:
  - All the eigenvalues are real and positive
  - All the **eigenvectors** are **orthogonal** to each other



## Mathematical background

- Starting from the definition  $\mathbf{X} = \mathbf{T} \cdot \mathbf{P}^T$ , one can obtain the following relationships

$$\mathbf{C} = \mathbf{X}^T \cdot \mathbf{X} = \mathbf{P} \cdot \mathbf{T}^T \cdot \mathbf{T} \cdot \mathbf{P}^T = \mathbf{P} \cdot \mathbf{\Lambda} \cdot \mathbf{P}^T = \mathbf{P} \cdot \mathbf{\Lambda} \cdot \mathbf{P}^{-1}$$

- The latter equation corresponds to the **eigendecomposition** of the square matrix  $\mathbf{C} = \mathbf{X}^T \cdot \mathbf{X}$ 
  - $\mathbf{\Lambda}$  is a diagonal matrix whose diagonal elements are the eigenvalues of **C**
  - The  $m$ -th element  $\mathbf{t}_m^T \mathbf{t}_m = \mathbf{\Lambda}_{mm}$  is the **variance** explained by the  $m$ -th score
  - **P** is the  $n \times n$  square matrix whose  $m$ -th column is the eigenvector  $\mathbf{p}_m$  of



• it is a rotation matrix

## Mathematical background

- Once the eigenvectors  $\mathbf{p}_m$  are computed the corresponding scores can be derived

$$\mathbf{X} = \mathbf{T} \cdot \mathbf{P}^T \Rightarrow \mathbf{X} \cdot \mathbf{P} = \mathbf{T} \cdot \mathbf{P}^T \cdot \mathbf{P} \Rightarrow \mathbf{T} = \mathbf{X} \cdot \mathbf{P}$$

- In practice, the original variables are **projected** onto the orthogonal eigenspace defined by the eigenvectors/loadings

## Mathematical background

- The eigenvalues of the covariance matrix are related to the variance of the scores

$$\lambda_j = \mathbf{t}_j^T \mathbf{t}_j = \sum_{i=1}^I t_{i,j}^2$$

- Thus the j-th eigenvalue is the **dispersion captured by the j-th score**
- **The total variance in the original data set is preserved in the T matrix**

Sum of the variances of  
the original variables

=

Sum of the variances  
of the scores

## Mathematical background

- In summary, one ends up with two matrices

$$\mathbf{T}_{(I \times J)} = \left[ \begin{array}{c|c|c} \mathbf{t}_1 & \mathbf{t}_2 & \dots & \mathbf{t}_j \end{array} \right] \quad \mathbf{P}_{(J \times J)} = \left[ \begin{array}{c|c|c} \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_j \end{array} \right]$$

### Scores matrix

The  $j$ -th column represents an **independent variable** obtained by projecting the data onto the  $j$ -th eigenvector

Remind:

Sort the eigenvectors according to their eigenvalue size (that is, their variance)

### Loading

Each column is an eigenvector of the covariance matrix



## PCA – Dimension reduction

- The scores and loading matrices can be approximated by considering only the **first  $A$  principal components**

$$\mathbf{T}_{(I \times I)} = \left[ \begin{array}{c|c} \mathbf{T}_A & \tilde{\mathbf{T}}_A \end{array} \right]_{(I \times A) \quad (I \times (I-A))} \approx \mathbf{T}_A_{(I \times A)}$$

$$\mathbf{P}_{(J \times J)} = \left[ \begin{array}{c|c} \mathbf{P}_A & \tilde{\mathbf{P}}_A \end{array} \right]_{(J \times A) \quad (J \times (J-A))} \approx \mathbf{P}_A_{(J \times A)}$$

Information considered negligible

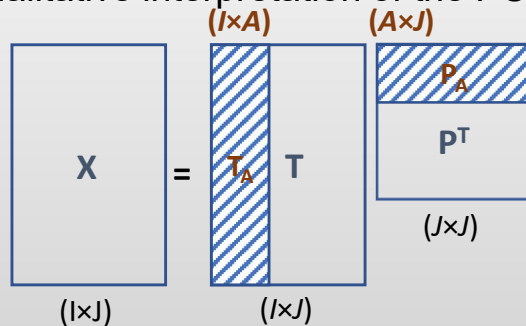
Information considered negligible





## PCA – Dimension reduction

- Qualitative interpretation of the PCA



$$\mathbf{X} \approx \mathbf{T}_A \mathbf{P}_A^T$$

$$(I \times J) \approx (I \times A)(A \times J)$$

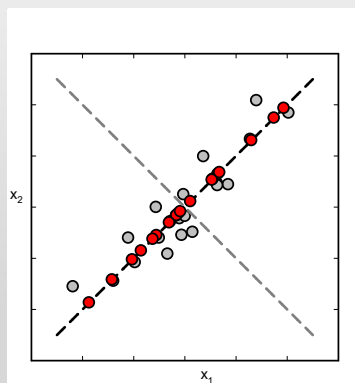
In general:  
 $A \ll J$

- Only part of the information collected in the  $\mathbf{X}$  matrix is relevant
- Only the first  $A$  columns of  $\mathbf{T}$  (the first scores) take into account most of the data variance



## PCA – A geometrical interpretation

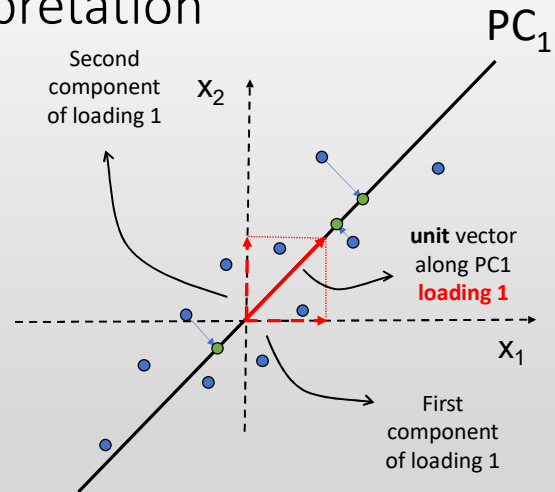
- 2D example - Reduction to 1D



- Samples are strongly correlated
- First principal component PC1 is the eigenvector direction corresponding to maximum variance (largest eigenvalue) in the coordinate space
- Second principal component is the orthogonal one leading to the second variance directions

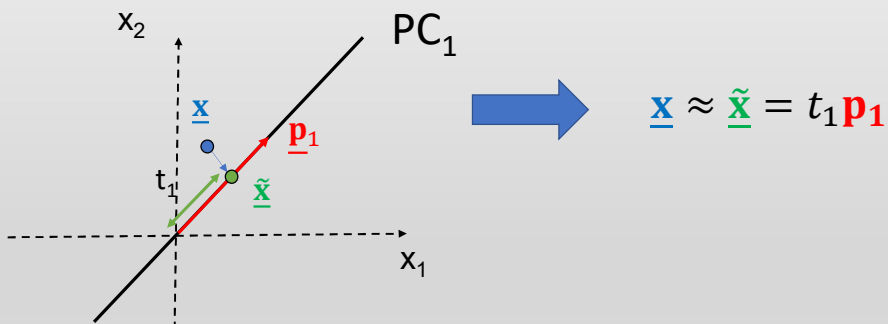
## PCA – A geometrical interpretation

- **Orthogonal projection** onto a specific PC results in a **score** for each sample
- The loading is the **unit vector** which defines this direction



## PCA – A geometrical interpretation

- The score is the projection of the point onto the **first** principal component



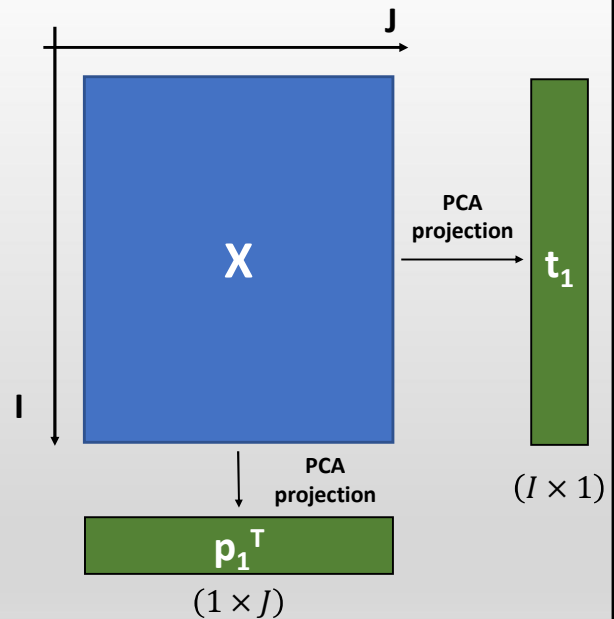
## PCA – Working principle – Reduction to 1D

- PCA projects matrix  $\mathbf{X}$  into:
  - a **score** vector  $\mathbf{t}_1$
  - a **loading** vector  $\mathbf{p}_1$

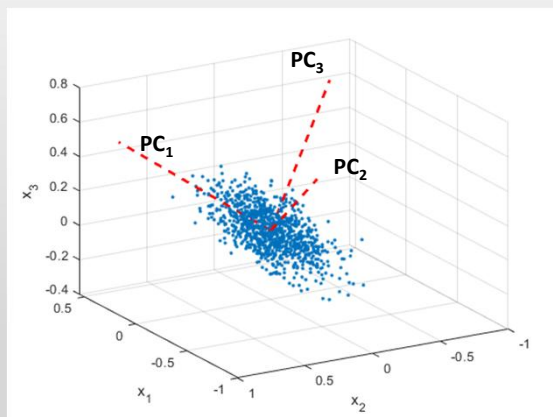
$$\underline{\mathbf{X}} \approx \underline{\mathbf{t}}_1 \underline{\mathbf{p}}_1^T$$

$(I \times J) \quad (I \times 1)(1 \times J)$

- $\mathbf{t}_1$  and  $\mathbf{p}_1$  are the first components



## PCA – A geometrical interpretation



- 3D example (A little bit more complicated)
- Points are mostly aligned along the 2D plane defined by the  $PC_1$  and the  $PC_2$  directions

## PCA – Working principle – Reduction to 2D

- If two principal components are required, matrix is formed by the outer products of  $\mathbf{t}_1$  and  $\mathbf{p}_1$ ,  $\mathbf{t}_2$  and  $\mathbf{p}_2$

$$\mathbf{X} = \begin{array}{c} \mathbf{p}_1 \\ \hline \mathbf{t}_1 \end{array} + \begin{array}{c} \mathbf{p}_2 \\ \hline \mathbf{t}_2 \end{array} + \mathbf{E}$$

- Matrix  $\mathbf{X}$  is decomposed into two sets of rank 1 outer products (2 terms) and the **residual matrix E**



## PCA – Working principle

- Successive components are formed by the outer products of  $\mathbf{t}_a$  and  $\mathbf{p}_a$

$$\mathbf{X} = \begin{array}{c} \mathbf{p}_1 \\ \hline \mathbf{t}_1 \end{array} + \begin{array}{c} \mathbf{p}_2 \\ \hline \mathbf{t}_2 \end{array} + \dots + \begin{array}{c} \mathbf{p}_A \\ \hline \mathbf{t}_A \end{array} + \mathbf{E}$$

- Matrix  $\mathbf{X}$  is decomposed into a set of rank 1 outer products (A terms) and the **residual matrix E**



## PCA – Working principle

- The master equation for PCA is eventually

$$\mathbf{X} = \mathbf{t}_1 \mathbf{p}_1^t + \mathbf{t}_2 \mathbf{p}_2^t + \dots + \mathbf{t}_A \mathbf{p}_A^t + \mathbf{E}$$

- or

$$\begin{array}{ccccccc} \mathbf{X} & = & \mathbf{T}_A & \mathbf{P}_A^T & + & \mathbf{E} & \\ (I \times J) & & (I \times A) & (A \times J) & & (I \times J) & \\ \swarrow & & \swarrow & \swarrow & & \swarrow & \\ \text{original data} & & \text{score} & \text{loading} & & \text{residual} & \\ \text{matrix} & & \text{matrix} & \text{matrix} & & \text{matrix} & \end{array}$$

## Estimation of the residuals

- When considering a PCA model with  $A$  principal components, one can evaluate the residual  $\mathbf{E}$

$$\mathbf{E} = \mathbf{X} - \mathbf{T}_A \cdot \mathbf{P}_A^T = \mathbf{X} \cdot (\mathbf{I} - \mathbf{P}_A \cdot \mathbf{P}_A^T)$$

## Estimation of the components

- How many principal components are needed?
- Possible criterion: **cumulative variance** explained by the **first A** principal components
  - The number of principal components to be considered explains most of the variance in the data (e.g., 95%)
- Alternative possibilities will be discussed in the case studies

## PCA to predict new data – Projection of the data onto the principal component space

- Single observations, (eventually new data  $\mathbf{x}_{new}$ ) can be eventually **projected** onto the space defined by the PCA model:

$$\mathbf{t}_{new} = \mathbf{x}_{new} \cdot \mathbf{P}_A$$

$(1 \times A) \quad (1 \times J) \quad (J \times A)$

$$\hat{\mathbf{x}}_{new} = \mathbf{x}_{new} \cdot \mathbf{P}_A \cdot \mathbf{P}_A^T$$

$(1 \times J) \quad (1 \times J) \quad (J \times A) \quad (A \times J)$

## PCA – Summary

- PCA projects the original data onto an orthogonal eigenspace of smaller dimensions
- The space is described by the first  $A$  eigenvectors of the covariance matrix
- The scores (i.e. the data projections onto the first eigenvectors) represent a set of independent variables
- New data can be projected in the PCA model

## References

1. Brereton, R.G. *Chemometrics: Data Analysis for the Laboratory and Chemical Plant*. Wiley, 2003
2. Brereton, R.G. *Chemometrics for Pattern Recognition*. Wiley, 2009
3. Jackson, J.E., *A User's Guide to Principal Components*. Wiley, New York, 1991
4. Jolliffe, I.T. *Principal Component Analysis*. Second Edition. Springer, 2002.
5. Jolliffe IT, Cadima J. (2016). Principal component analysis: a review and recent developments. *Phil. Trans.R.Soc. A374*:20150202.
6. Wold S., Esbensen K., Geladi P (1987). Principal Component Analysis – A tutorial. *Chemom. Intell. Lab. 2*, 37-52