

Principal Component Analysis

Machine Learning 2023-24

UML Book Chapter 23

Slides P. Zanuttigh (derived from F. Vandin slides)

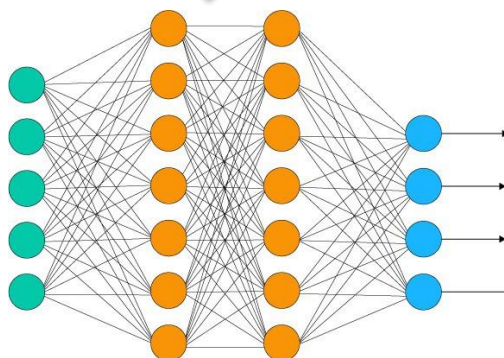
Recall: Supervised Learning



Training data
with labels



Training procedure



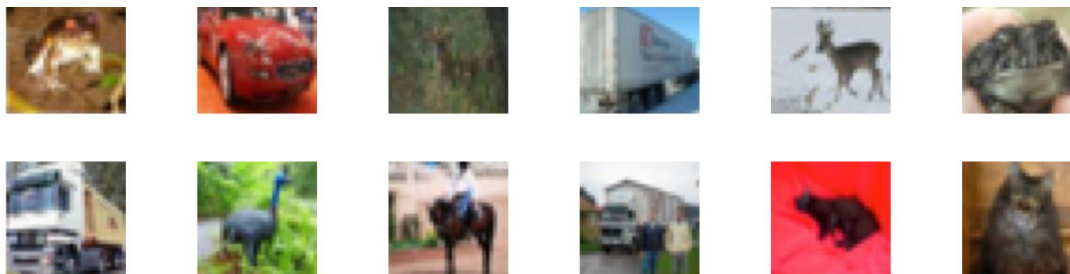
Data to be
analyzed

ML model

(training: estimate parameters)

Predicted
Label

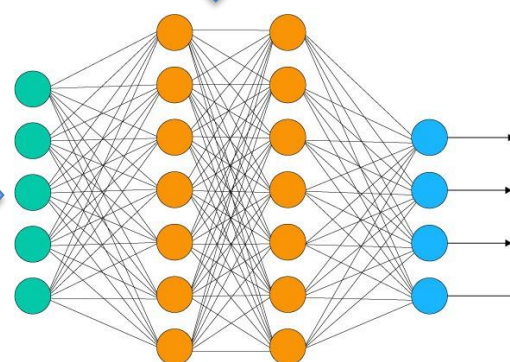
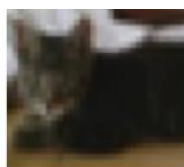
Unsupervised Learning: Training Data is Unlabeled



Training data
(**unlabeled**)



Training procedure



ML model



Prediction

Data to be
analyzed

(training: estimate parameters)

Unsupervised Learning:
Training data is not labeled

Unsupervised Learning Techniques

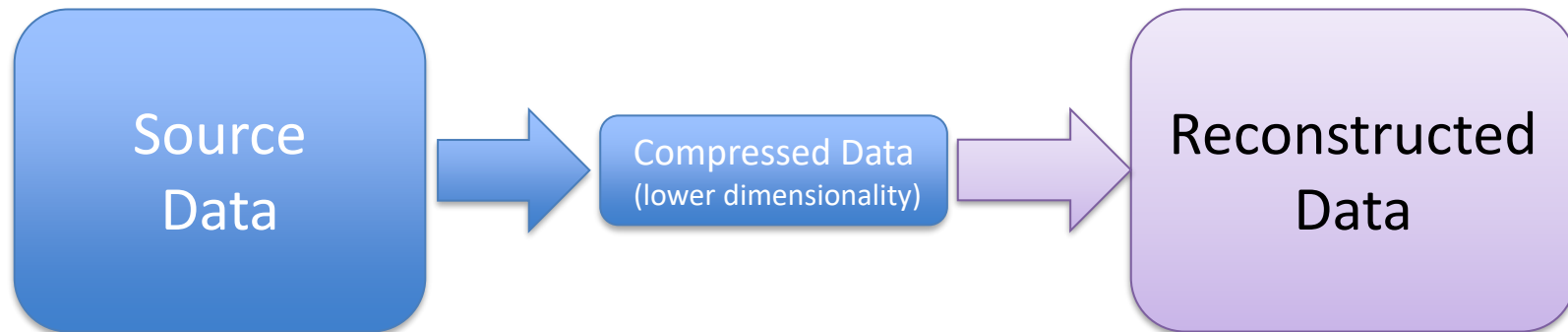
We are going to see only a couple of unsupervised learning tasks and a few very simple and commonly used methods

- *Clustering (already seen)*
 - *K-means*
 - *Linkage-based clustering*
- *Dimensionality reduction*
 - *Principal Component Analysis (PCA)*

There are many other techniques (*not part of this course*)

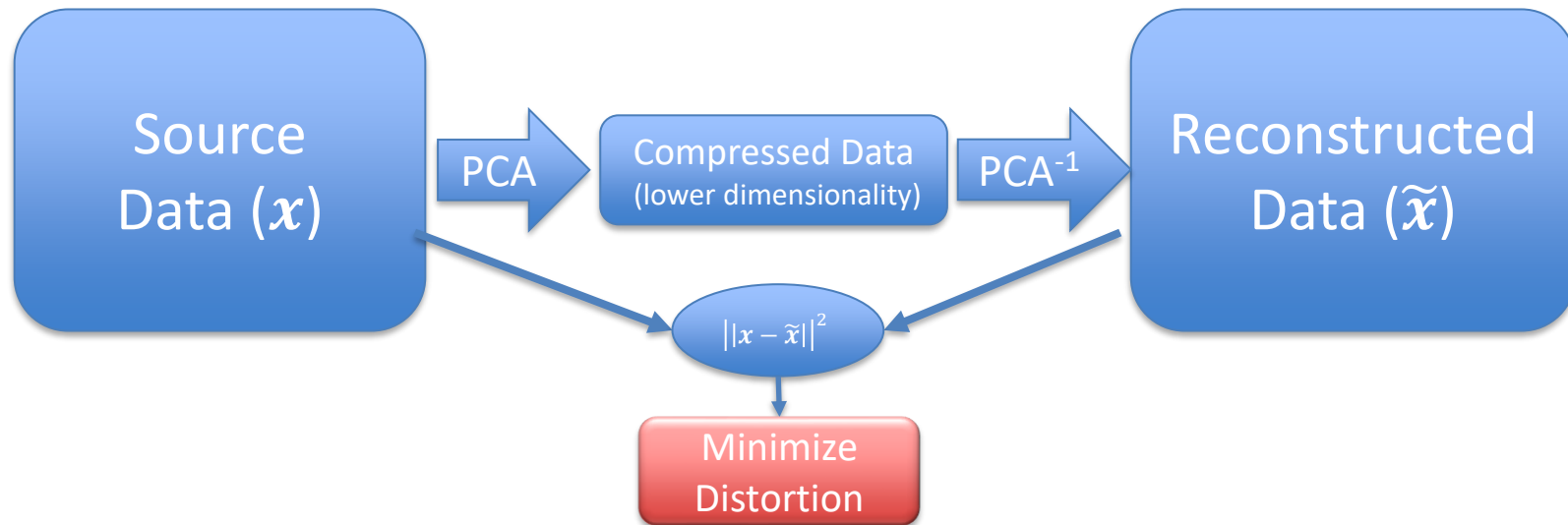
- Mean shift clustering, spectral clustering....
- Compressive sensing

Dimensionality Reduction



- ❑ Take data from a **highly** dimensional space and project to a **lower** dimensional one
- ❑ Many applications:
 - Reduce number of features (learn with less samples, lower computation req.)
 - Capture most important aspects of the data for subsequent analysis
 - Data visualization
 - etc..
- ❑ By reducing dimensionality part of the information get lost
 - Lower dimensional data should be a **good approximation** of the higher dimensional representations
 - *Good approximation: lower dimensional data should allow for reconstructing the original data with a reasonable accuracy*

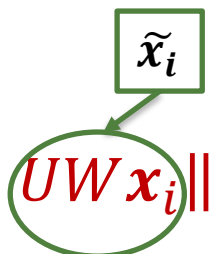
Dimensionality Reduction



- ❑ Lower dimensional data should be a *good approximation* of the higher dimensional representations
- ❑ *Good approximation*: minimize error obtained by reprojecting the data back to the high dimensional space (→ similar to lossy data compression)
- ❑ Focus on *linear mapping* of the data (represented by a matrix multiplication)
- ❑ *Principal Component Analysis (PCA)*: find the linear mapping that minimizes the mean squared error in the reprojection

Principal Component Analysis (PCA)

- ❑ $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m \in \mathbb{R}^d$: data points
- ❑ $W \in \mathbb{R}^{n,d}$ ($n < d$): mapping $\mathbf{x} \rightarrow \mathbf{y} = W\mathbf{x}$
 - where $\mathbf{y} = W\mathbf{x} \in \mathbb{R}^n$ is a lower dimensional representation of $\mathbf{x} \in \mathbb{R}^d$
- ❑ $U \in \mathbb{R}^{d,n}$ ($n < d$): inverse mapping $\mathbf{y} \rightarrow \tilde{\mathbf{x}} = U\mathbf{y}$
 - used to recover an approximation $\tilde{\mathbf{x}} = UW\mathbf{x}$ of \mathbf{x}
- ❑ Target: find the lower dimensional representation that better approximates the data \rightarrow that leads to minimum squared distance between $\tilde{\mathbf{x}}$ and \mathbf{x}
 - Corresponds to seek for the n -dimensional basis that best captures the variance in the d -dimensional data

$$\operatorname{argmin}_{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}} \sum_{i=1}^m \|\mathbf{x}_i - UW\mathbf{x}_i\|_2^2$$


Lemma

There exist an optimal solution (U^*, W^*) of $\operatorname{argmin}_{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}} \sum_{i=1}^m \|x_i - UWx_i\|_2^2$ where:

Recall: orthonormal

- the columns of U^* are orthonormal (i.e., $(U^*)^T U^* = I$)
- $W^* = (U^*)^T$

- $u_i^T u_j = 0, \forall i \neq j$
- $\|u_i\| = 1 = u_i^T u_i, \forall i$

Demonstration:

1. Fix U, W and consider the mapping $x \rightarrow UWx$
 - The range of the mapping is $R = \{UWx ; x \in \mathbb{R}^d\}$
2. $V \in \mathbb{R}^{d,n}$: matrix whose column form an orthonormal basis of R
 - Recall that $V^T V = I$ and $\forall x \in R: \exists y \in \mathbb{R}^n$ with $x = Vy$
3. $\forall x \in \mathbb{R}^d, \forall y \in \mathbb{R}^n$:
 - $\|x - Vy\|_2^2 = \|x\|^2 + y^T V^T V y - 2y^T V^T x = \|x\|^2 + \|y\|^2 - 2y^T (V^T x)$
4. Minimize $\|x\|^2 + \|y\|^2 - 2y^T (V^T x)$ w.r.t y : set $\nabla = 0 \rightarrow 2y - 2(V^T x) = 0 \rightarrow y_{opt} = V^T x$
5. $\forall x: \operatorname{argmin}_{\tilde{x} \in R} \|x - \tilde{x}\|_2^2 = Vy_{opt} = V(V^T x)$: it is the best approximation in subspace R
6. $\forall x$: includes also x_1, \dots, x_m (data vectors): $\sum_{i=1}^m \|x_i - UWx_i\|^2 \geq \sum_{i=1}^m \|x_i - VV^T x_i\|^2$, so we can replace U, W with VV^T without increasing the objective
7. Holds for $\forall U, W$: there exist a solution that minimize $\sum_{i=1}^m \|x_i - UWx_i\|^2$ with V orthonormal columns and $W = U^T$

Optimization Problem

There exist an optimal solution (U^*, W^*) of $\operatorname{argmin}_{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}} \sum_{i=1}^m \|x_i - UWx_i\|_2^2$ where:

- the columns of U^* are orthonormal (i.e., $(U^*)^T U^* = I$)
- $W^* = (U^*)^T$

The optimization problem can be rewritten as:

$$\operatorname{argmin}_{U \in \mathbb{R}^{d,n}: U^T U = I} \sum_{i=1}^m \|x_i - UU^T x_i\|_2^2$$

- Trace: Σ elements on diagonal
- It is a scalar
- $\operatorname{trace}(A^T B) = \operatorname{trace}(AB^T) = \operatorname{trace}(B^T A) = \operatorname{trace}(BA^T)$

With some manipulations: $\|x - UU^T x\|_2^2 = \|x\|^2 - 2x^T UU^T x + x^T UU^T UU^T x = \|x\|^2 - x^T UU^T x = \|x\|^2 - \operatorname{trace}(x^T UU^T x) = \|x\|^2 - \operatorname{trace}(U^T x x^T U)$

$$\operatorname{argmax}_{U \in \mathbb{R}^{d,n}: U^T U = I} \operatorname{trace} \left(U^T \sum_{i=1}^m x_i x_i^T U \right) = \operatorname{argmax}_{U \in \mathbb{R}^{d,n}: U^T U = I} \operatorname{trace}(U^T A U)$$

Notice: $A = \sum_{i=1}^m x_i x_i^T$ is symmetric and positive semidefinite. It can be rewritten as $A = V D V^T$ where D is diagonal (with eigenvalues $D_{d,d} \geq 0$) and $V^T V = V V^T = I$ (the columns of V are the eigenvectors of A)

Theorem (PCA)

Let $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ be arbitrary vectors in \mathbb{R}^d

let $A = \sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i^T$

let $\mathbf{u}_1, \dots, \mathbf{u}_n$ be n eigenvectors of A corresponding to the largest n eigenvalues of A



Then a solution of the PCA optimization $\underset{W \in \mathbb{R}^{n,d}, U \in \mathbb{R}^{d,n}}{\operatorname{argmin}} \sum_{i=1}^m \|\mathbf{x}_i - UW\mathbf{x}_i\|_2^2 = \underset{U \in \mathbb{R}^{d,n}; U^T U = I}{\operatorname{argmax}} \operatorname{trace}(U^T A U)$
is to set U to be the matrix whose columns are $\mathbf{u}_1, \dots, \mathbf{u}_n$ and to set $W = U^T$

Notes:

- ❑ Recall: Decompose A as VDV^T (SVD decomposition, D diag. and $V^T V = VV^T = I$)
- ❑ It is a common practice to "center" the examples before applying PCA (i.e., subtract the mean)
- ❑ Computation time is $O(d^3) + O(md^2)$ (the first term for calculating eigenvalues and the second for constructing A)
- ❑ Trick for faster solution in case $d \gg m$ (not part of the course)



Pseudocode

Input

$X \in \mathbb{R}^{m,d}$: matrix that contains m samples, one for each row

n : number of components

Algorithm

Compute $A = X^T X$

Perform eigenvalue decomposition of A

Let u_1, \dots, u_n be the eigenvectors of A corresponding to largest eigenvalues

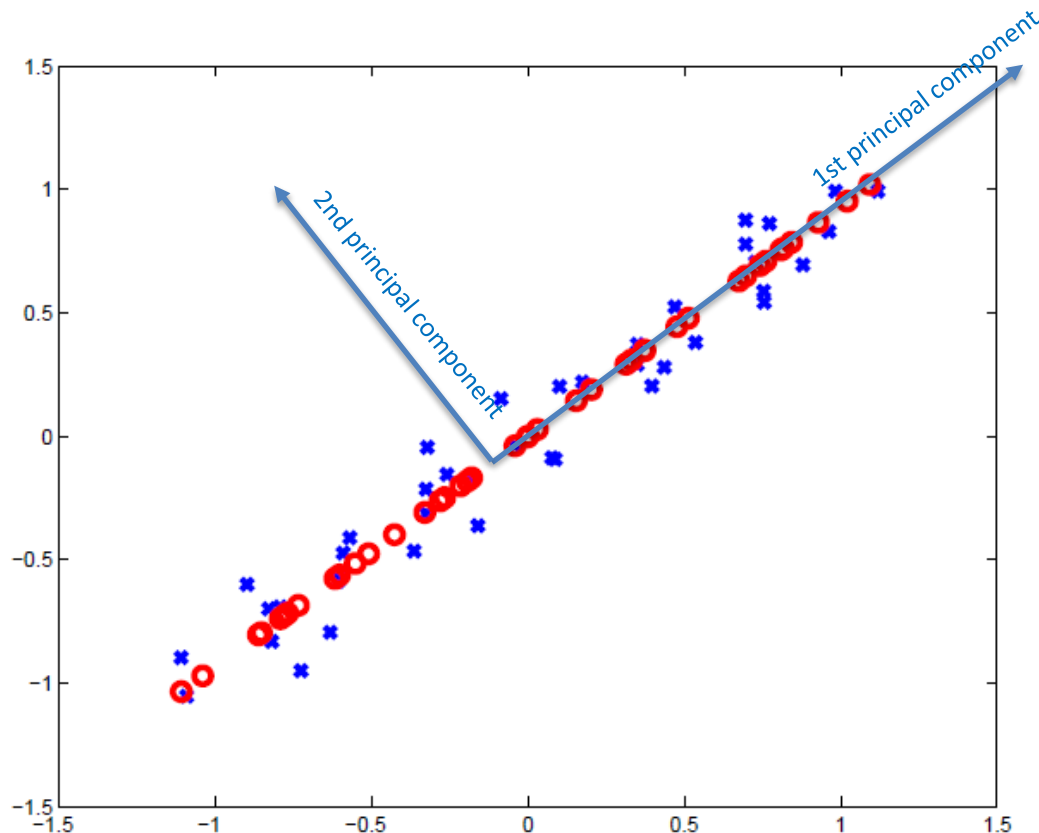
Output: u_1, \dots, u_n

On the book: trick for $d \gg m$
(not part of the course)

Eigenvectors of A :

1. First principal component (p.c.) = direction with largest projected variance
2. Second p.c. = orthogonal direction with largest projected variance
 - i.e., largest remaining variance after removing the first p.c.
3. ... iterate for all the other components (3...n)

Example: From 2D to 1D

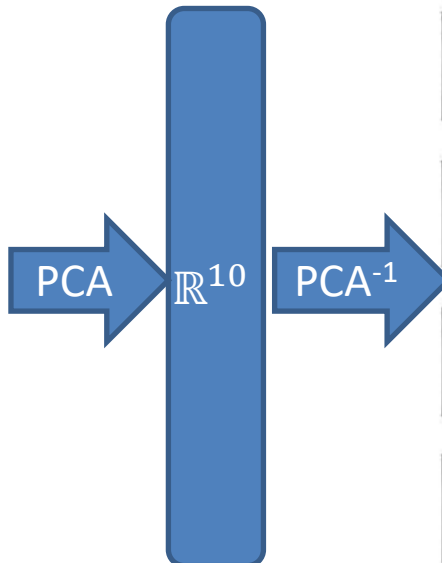


*Set of 2D vectors (**blue**) and their reconstruction (**red**)
after dimensionality reduction to 1D with PCA*

Example: Face Compression



$$\mathbb{R}^{50 \times 50} = \mathbb{R}^{2500}$$



$$\mathbb{R}^{50 \times 50} = \mathbb{R}^{2500}$$

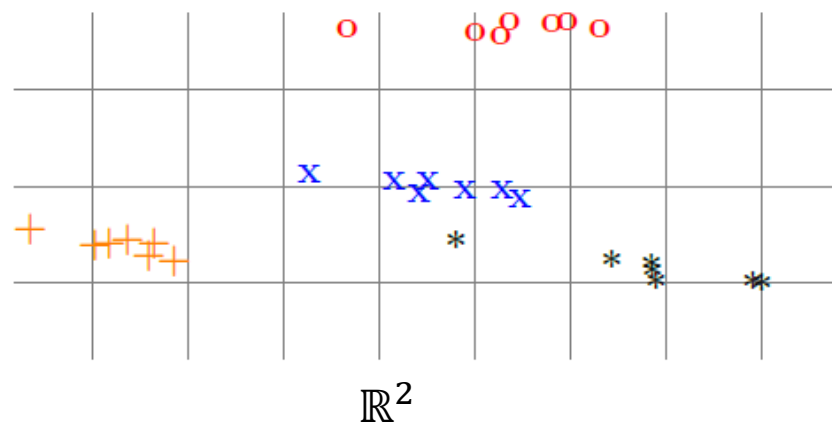
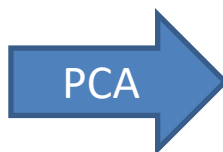
*enlarged
example*



Example: Face Recognition



$$\mathbb{R}^{50 \times 50} = \mathbb{R}^{2500}$$



- Faces with the same type of mark (+, x, *, o) belong to the same individual
- PCA can be used for face recognition ! (*eigenfaces* algorithm)