



Aprendizagem 2021

Lab 10: Dimensionality Reduction

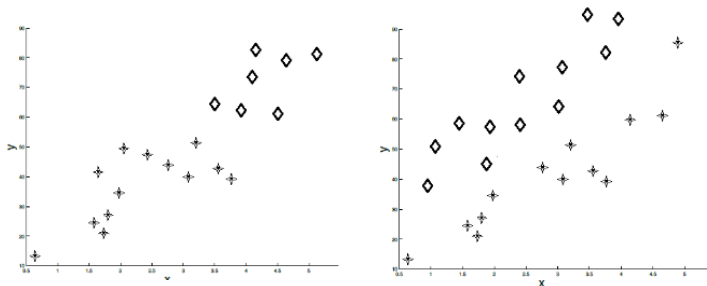
Practical exercises

1. Given the following training data

	y_1	y_2
x_1	0	0
x_2	4	0
x_3	2	1
x_4	6	3

- Compute the K-L transformation
- What is the rotation applied to go from the original to the eigenvector coordinate system?
- Which eigenvector is most significant?
- Can we apply the Kaiser criterion?
- Map the points onto the most significant dimension

2. Given the following datasets where observations are in \mathbb{R}^2 and belong to one of two classes:



Which principal components can accurately discriminate the class per dataset?

3. The following top-7 eigenvalues explain 90% of the variation of dataset X :

$$\lambda_1=20, \lambda_2=10, \lambda_3=5, \lambda_4=4, \lambda_5=3, \lambda_6=2, \lambda_7=1$$

What is the most accurate information regarding X :

- X has less than 7 attributes
- X has 7 attributes
- X has more than 7 attributes
- X has more than 11 attributes

4. Given a set of data points in \mathbb{R}^3 , the following covariance matrix was obtained:

$$\begin{bmatrix} 91.43 & 171.92 & 297.99 \\ & 373.92 & 545.21 \\ & & 1297.26 \end{bmatrix}$$

as well as the following eigenvectors retrieved:

$$\mathbf{u}_1 = \begin{pmatrix} 0.2179 \\ 0.4145 \\ 0.8836 \end{pmatrix}, \mathbf{u}_2 = \begin{pmatrix} -0.2466 \\ -0.8525 \\ 0.4608 \end{pmatrix}, \mathbf{u}_3 = \begin{pmatrix} 0.9443 \\ -0.3183 \\ -0.0836 \end{pmatrix}$$

Please select the more complete answer:

- i. eigenvalue λ_1 is approximately 1626
- ii. eigenvalue λ_2 is approximately 129
- iii. eigenvalues λ_1 and λ_2 explain >99% of the variation in data
- iv. all of the above

5. Given the following dataset:

	y_1	y_2
x_1	1	-1
x_2	0	1
x_3	-1	0

and the corresponding eigenvectors and eigenvalues:

$$\lambda_1 = 3/2 \text{ and } \lambda_2 = 1/2$$

$$\mathbf{u}_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \mathbf{u}_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

- a) Transform the input data using PCA
- b) [optional] Assess the data recovery error when considering the most informative component only

Programming quest

6. Considering the *housing* dataset available at <https://web.ist.utl.pt/~rmch/dscience/data/housing.arff>
 - a. How many principal components are necessary to explain 90% of data's variability?
 - b. Compare the accuracy of one of the covered classifiers in the original *versus* reduced data space. Why we are unable to observe considerable improvements?

Resource: <https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.PCA.html>